



Information Systems Research

Publication details, including instructions for authors and subscription information:
<http://pubsonline.informs.org>

Dynamic, Multidimensional, and Skillset-Specific Reputation Systems for Online Work

Marios Kokkodis

To cite this article:

Marios Kokkodis (2021) Dynamic, Multidimensional, and Skillset-Specific Reputation Systems for Online Work. Information Systems Research 32(3):688-712. <https://doi.org/10.1287/isre.2020.0972>

Full terms and conditions of use: <https://pubsonline.informs.org/Publications/Librarians-Portal/PubsOnLine-Terms-and-Conditions>

This article may be used only for the purposes of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval, unless otherwise noted. For more information, contact permissions@informs.org.

The Publisher does not warrant or guarantee the article's accuracy, completeness, merchantability, fitness for a particular purpose, or non-infringement. Descriptions of, or references to, products or publications, or inclusion of an advertisement in this article, neither constitutes nor implies a guarantee, endorsement, or support of claims made of that product, publication, or service.

Copyright © 2021, INFORMS

Please scroll down for article—it is on subsequent pages



With 12,500 members from nearly 90 countries, INFORMS is the largest international association of operations research (O.R.) and analytics professionals and students. INFORMS provides unique networking and learning opportunities for individual professionals, and organizations of all types and sizes, to better understand and use O.R. and analytics tools and methods to transform strategic visions and achieve better outcomes.

For more information on INFORMS, its publications, membership, or meetings visit <http://www.informs.org>

Dynamic, Multidimensional, and Skillset-Specific Reputation Systems for Online Work

Marios Kokkodis^a

^a Carroll School of Management, Boston College, Chestnut Hill, Massachusetts 02467

Contact: kokkodis@bc.edu,  <https://orcid.org/0000-0002-5037-6060> (MK)

Received: January 15, 2019

Revised: December 17, 2019; May 31, 2020

Accepted: August 4, 2020

Published Online in Articles in Advance:
March 31, 2021

<https://doi.org/10.1287/isre.2020.0972>

Copyright: © 2021 INFORMS

Abstract. Reputation systems in digital workplaces increase transaction efficiency by building trust and reducing information asymmetry. These systems, however, do not yet capture the dynamic multidimensional nature of online work. By uniformly aggregating reputation scores across worker skills, they ignore skillset-specific heterogeneity (reputation attribution), and they implicitly assume that a worker's quality does not change over time (reputation staticity). Even further, reputation scores tend to be overly positive (reputation inflation), and, as a result, they often fail to differentiate workers efficiently. This work presents a new augmented intelligence reputation framework that combines human input with machine learning to provide dynamic, multidimensional, and skillset-specific worker reputation. The framework includes three components: The first component maps skillsets into a latent space of finite competency dimensions (word embedding), and, as a result, it directly addresses reputation attribution. The second builds dynamic competency-specific quality assessment models (hidden Markov models) that solve reputation staticity. The final component aggregates these competency-specific assessments to generate skillset-specific reputation scores. Application of this framework on a data set of 58,459 completed tasks from a major online labor market shows that, compared to alternative reputation systems, the proposed approach (1) yields more appropriate rankings of workers that form a closer-to-normal reputation distribution, (2) better identifies "nonperfect" workers who are more likely to underperform and are harder to predict, and (3) improves the ranking of within-opening choices and yields significantly better outcomes. Additional analysis of 77,044 restaurant reviews shows that the proposed framework successfully generalizes to alternative contexts, where assigned feedback scores are overly positive and service quality is multidimensional and dynamic.

History: This article was accepted by Special Section Editors Hemant Jain, Balaji Padmanabhan, Paul A. Pavlou, and Raghu T. Santanam for the *Information Systems Research* Special Section on Humans, Algorithms, and Augmented Intelligence: The Future of Work, Organizations, and Society.

Supplemental Material: The online appendices are available at <https://doi.org/10.1287/isre.2020.0972>.

Keywords: reputation frameworks • reputation inflation • reputation attribution • reputation staticity • online labor markets • hidden Markov models • word embedding

1. Introduction

Online labor markets (Peopleperhour, Freelancer) facilitate global short-term contracts or freelance work (Graham et al. 2017). Buyers purchase services from an abundance of capable online workers who complete diverse tasks, including web development, graphic design, accounting, sales, marketing, and data science. On par with other online platforms, online labor marketplaces grew exponentially during the past decade (Upwork 2014, Freelancers Union 2017). This growth will likely continue (if not accelerate) in the future, as automation and the sharing economy structure the future of work (Sundararajan 2016, IBM Institute of Business Value 2019).

One determinant of success for online labor markets is the intermediary trust that they instill between employers and workers (Ba and Pavlou 2002, Pavlou

and Gefen 2004, Nica et al. 2017). Reputation systems are a standard mechanism that online labor markets use to increase trust and reduce information asymmetry (Akerlof 1970, Kokkodis and Ipeirotis 2016, Yoganarasimhan 2013, Nica et al. 2017). These systems rely on human input: employers rate workers for the tasks they complete, and these ratings become part of the workers' online resumes. Such reputation mechanisms measure expected service quality (Kokkodis and Ipeirotis 2016, Rahman 2018b, Filippas et al. 2018), and, as a result, they increase employers' trust in workers' abilities and facilitate market transactions (Yoganarasimhan 2013, Moreno and Terwiesch 2014, Lin et al. 2016). Besides, workers realize how reputation instills trust and affects employer choices, and they tend to readjust their premiums according to their reputation scores (Banker and Hwang 2008, Gandini et al. 2016).

Despite these benefits, the design of current reputation systems does not capture the dynamic and multidimensional nature of online work. In particular, reputation systems in online labor markets implicitly assume (by uniformly averaging all prior feedback scores) that worker quality and expertise are not evolving (Hendrikx et al. 2015). However, current trends show that new skills are born and old skills die faster than ever (Autor et al. 1998, Autor 2001, Oliver 2015, Kokkodis and Ipeirotis 2016). As a result, to remain marketable, online workers must be diligently and continuously re-educating and reskilling themselves (Kuhn and Skuterud 2004, Stevenson 2009, Oliver 2015). Furthermore, online worker reputation scores are unidimensional and skillset-independent. Digital workplaces, however, are highly heterogeneous in terms of qualifications (Kokkodis and Ipeirotis 2014), whereas online workers tend to complete tasks that require diverse skillsets (Kokkodis and Ipeirotis 2016). As a result, these unidimensional reputation scores cannot provide accurate estimates of skillset-specific expertise. Finally, on par with other online platforms (Hu et al. 2009, Hu et al. 2017, Zervas et al. 2015), online worker reputation scores tend to be overly positive (Abhinav et al. 2017, Filippas et al. 2018). Such inflated scores do not sufficiently differentiate workers, as most of them are rated as “better than average” (Filippas et al. 2018).

Given these shortcomings of current reputation systems in online labor markets, *how can we design dynamic, multidimensional, skillset-specific reputation frameworks?* To address this question, I propose an intelligence augmentation (IA) system that relies on three design principles: (1) mapping of any combination of arbitrary skills into a latent space of finite competency dimensions (word embedding), (2) dynamic competency-specific quality assessment (hidden Markov models), and (3) aggregation of these competency-specific assessments. Decomposition of skills to competencies facilitates skillset-specific reputation. Dynamic quality assessment explains the evolution of workers as they learn new skills and gain expertise. Aggregation of the competency-specific quality assessments for any given combination of skills results in representative (normal-like) skillset-specific reputation distributions (Schmidt and Hunter 1983) that catalyzes worker differentiation.

Analysis of 58,459 completed tasks from a major online labor market shows that the proposed approach significantly outperforms 10 alternative advanced reputation systems (including the market’s current reputation system, systems that rely on link analysis, gradient boosting, neural networks, and adaptations of recommender systems). In particular, compared with these systems, the proposed approach (1) yields more appropriate rankings of workers that

form a closer-to-normal reputation distribution, (2) better identifies “nonperfect” workers who are more likely to underperform and are harder to predict, and (3) improves the ranking of within-opening choices and yields significantly better outcomes. Additional analysis of 77,044 restaurant reviews shows that the proposed framework successfully generalizes to alternative contexts, where assigned feedback scores are overly positive and service quality is multidimensional and dynamic.

This work is the first to identify shortcomings of current reputation systems of online labor markets and to present design principles that future reputation systems should have in order to estimate a worker’s dynamic and multidimensional reputation. By solidifying these principles into different components, the proposed IA framework combines human input with machine intelligence to result in accurate, skillset-specific reputation scores. Such accurate scores (1) help workers to differentiate, (2) guide employers to make informed and fast (reduced search cost; see Bakos 1997) decisions, and (3) enable the market to improve its recommendation algorithms and also understand the supply distributions across latent competencies. By predicting underperforming workers, the framework preemptively informs employers, an intervention that could reduce the number of adverse outcomes. Positive outcomes increase participation in the marketplace, thereby generating a continuous stream of revenue for the platform (Tripp and Grégoire 2011).

This IA framework also highlights how combining human input with advanced machine learning techniques can augment intelligence by creating the necessary conditions for humans to make informed decisions. Such systems have the potential to increase efficiency and outcome quality precisely because they intelligently differentiate workers (i.e., identify each individual’s latent qualities). Efficient differentiation can further guide labor supply redistribution (e.g., by motivating workers to re-educate) and inform career path advisers (Kokkodis and Ipeirotis 2020). As a result, the deployment of the proposed IA framework in different types of online platforms could have implications for workers, employers, businesses, and the future of work.

2. Research Context

Digital markets increase trust and signal the quality of their services and products through online reputation systems (Ba and Pavlou 2002, Dellarocas 2003, Pavlou and Gefen 2004, Dellarocas 2006, Zervas et al. 2015, Tadelis 2016). These systems facilitate product selection across a wide range of domains, including movies (Duan et al. 2008), books (Chevalier and Mayzlin 2006), music (Kokkodis and Ransbotham 2020), electronics (Ghose and Ipeirotis 2011, Cui et al. 2012),

hotels (Ye et al. 2009), local businesses (Lu et al. 2013, Luca 2016), and mobile apps (Lee and Raghu 2014).

2.1. Overview of Current Reputation Systems Designs

Given this established impact of online reputation systems, prior research has focused on improving their design and increasing their performance. Researchers have proposed reputation systems that have context-specific objectives and serve alternative domains such as e-commerce platforms, online communities, crowdsourcing platforms, and peer-to-peer networks. Based on their architecture, I cluster existing reputation systems into *human-based*, *machine-based*, and *hybrid* (human and machine).

2.1.1. Human-Based Reputation Systems. Most commercial reputation systems rely solely on human ratings (e.g., in online marketplaces; see Einav et al. 2016, Tadelis 2016, Luca 2017). For instance, Amazon users post reviews, rate products, and rate other reviews (Amazon 2018). Similarly, eBay sellers and buyers leave feedback for each other (eBay 2018). This feedback reflects both an overall rating (good, neutral, and negative) and numerical ratings for accuracy, communication, shipping time, and shipping charges. A similar reputation system appears in many online question-and-answer communities (e.g., Stack Overflow and discourse communities), where users up-vote or down-vote responses (Stack Overflow 2018, Kokkodis et al. 2020b). Third-party reputation platforms that allow users to review and rate are also available for multiple products and services, such as restaurants and hotels (e.g., TripAdvisor, Yelp; see Kokkodis and Lappas 2020) and Amazon Mechanical Turk (AMT) users (Turkopticon 2018).

2.1.2. Machine-Based Reputation Systems. Machine-based approaches do not require human raters. Instead, they often rely on network analysis to identify user quality. In online and question-answering communities, such methods focused on identifying expert (or helpful) users (Jurczyk and Agichtein 2007, Zhang et al. 2007, Bouguessa et al. 2008). In large organizations, proposed approaches combine information retrieval and graph-based techniques to analyze user social (and other) profiles and identify areas of expertise (Balog and De Rijke 2007). Similarly, peer-to-peer network reputation systems use network analysis to estimate the sharing quality of each participating node (Kamvar et al. 2003).

2.1.3. Hybrid Reputation Systems. Many reputation systems combine information from human raters with machine learning and network analysis. Such hybrid approaches are examples of IA systems (Jain et al. 2018),

as they enhance human judgment by combining artificial and human intelligence. The most basic ones combine ratings with information from social and other online sources (Sabater and Sierra 2001a,b; Hendriks and Bubendorfer 2013). These systems are tailored for e-commerce platforms (Hendriks and Bubendorfer 2013). For online communities, hybrid approaches use human cognitive traits along with subjective logic to identify experts (Pelechrinis et al. 2015). At the same time, peer-to-peer networks require different types of hybrid reputation systems that combine network analysis along with ratings and the personal histories of each node to estimate node trustworthiness (Damiani et al. 2002, Curtis et al. 2004, Xiong and Liu 2004, Tian and Yang 2011).

Hybrid reputation systems also appear in crowdsourcing settings (e.g., AMT; see Allahbakhsh et al. 2012, Jagabathula et al. 2014). For instance, some reputation management models include a rater's credibility along with information from each worker's set of completed tasks to estimate worker quality (Allahbakhsh et al. 2012). Similarly, and in order to filter out adversarial workers (Jagabathula et al. 2014), proposed reputation systems in crowdsourcing settings penalize workers with a poor reputation (Xie et al. 2015).

2.2. Reputation Systems in Online Labor Markets

Similar to most e-commerce platforms, online labor markets offer reputation systems that allow workers to receive feedback for the tasks they complete (Filippas et al. 2018, Wood-Doughty 2018). Over a series of completed tasks, these feedback scores accumulate to generate a worker's reputation on the platform (Rahman 2018a). Worker reputation "institutes trust among quasi-strangers" (Nica et al. 2017, p. 64), and, as a result, increases marketplace efficiency by reducing information asymmetry (Kokkodis et al. 2015). In particular, worker reputation is a major driving force in hiring choices (Yoganarasimhan 2013), and it correlates positively with worker earnings (Banker and Hwang 2008, Moreno and Terwiesch 2014, Gandini et al. 2016). Even having a reputation (compared to being new in the market) improves a worker's current (Lin et al. 2016) and subsequent hiring chances (Pallais 2014). These reputation effects are not uniform, as positive verified information appears to disproportionately benefit workers from less developed countries (Agrawal et al. 2013, Kanat et al. 2018). Finally, reputation is not necessarily category-specific; it transfers across multiple task categories that require diverse skill sets (Kokkodis and Ipeirotis 2016).

2.2.1. Shortcomings of Reputation Systems in Online Labor Markets. Despite these multidimensional effects, reputation systems in online labor markets are

not perfect (Filippas et al. 2018), as they experience: (1) reputation inflation, (2) reputation attribution, and (3) reputation staticity.

Reputation Inflation. Similar to other online marketplaces (Zervas et al. 2015, Hu et al. 2017), reputation scores in online labor markets are highly inflated (Abhinav et al. 2017, Filippas et al. 2018). This inflation happens mainly for two reasons. First, users who receive low feedback scores cannot get hired, and, as a result, they abandon the marketplace (Jøsang et al. 2007, Jøsang and Golbeck 2009, Jerath et al. 2011). Second, employers feel peer pressure to assign positive ratings (Filippas et al. 2018). The combination of the two yields reputation distributions that are positively skewed, where every worker is assumed to be “better than average.” Such inflated reputation scores do not sufficiently differentiate workers, as they form noisy estimates of service quality (Hendrikx et al. 2015).

Reputation Attribution. At the same time, current reputation systems in online labor markets provide unidimensional reputation scores that describe overall service quality. However, these digital workplaces are highly heterogeneous in terms of qualifications (Kokkodis and Ipeirotis 2014), as they offer tasks that require a diverse range of skills (e.g., logo design, software development, data analytics, marketing skills). Besides, workers often do not focus on specific types of tasks; instead, they complete tasks that require diverse skillsets (Kokkodis and Ipeirotis 2016). The existence of this highly heterogeneous environment in terms of skills suggests that unidimensional reputation scores cannot capture skill-specific qualities. For instance, consider a worker who provides an IT service and completes a task that requires networking, C, and Python. When the task is over, this worker receives a feedback score of 0.9. Does this score capture the worker’s service quality on networking, on C, on Python, or on any combination of these skills?

Reputation Staticity. Finally, the rapid evolution of skills and worker expertise in online labor markets further limits current reputation systems. Because new skills are born and old skills die faster than ever before (Autor et al. 1998, Autor 2001, Oliver 2015, Kokkodis and Ipeirotis 2016), workers need to continuously keep re-educating themselves (Kuhn and Skuterud 2004, Stevenson 2009, Oliver 2015, Kokkodis 2020). Workers are therefore dynamic entities that evolve by either gaining expertise on skills they know, or by investing in learning new skills (Kokkodis and Ipeirotis 2020). Current reputation systems assume that the quality of a service does not change over time (Jøsang et al. 2007), as they uniformly average received

ratings to provide an aggregate quality score (Hendrikx et al. 2015). This assumption is valid for a product in an e-commerce platform (e.g., a book or a camera), but it is misleading in representing the quality of a worker who gains expertise and acquires new skills over time.

These shortcomings of reputation systems in online labor markets result in service quality estimates that are often not predictive of future worker performance (Filippas et al. 2018). Consequently, decisions based on such reputation scores could yield unsuccessful collaborations that hurt the marketplace (Tripp and Grégoire 2011). As a result, there is a need for exploring alternative reputation systems that could potentially address these shortcomings and provide more representative reputation scores.

2.2.2. Do Existing Reputation Systems Address These Shortcomings?

Section 2.1 classifies current reputation systems into human-based, machine-based, and hybrid. The presented commercial applications of *human-based* reputation systems experience reputation inflation and, to a certain degree, reputation attribution and reputation staticity. In particular, e-commerce reputation scores are inflated (Chevalier and Mayzlin 2006, Hu et al. 2009, Zervas et al. 2015, Hu et al. 2017) due to response bias—that is, who chooses to rate a service (Moe and Schweidel 2012)—and due to acquisition bias—that is, buyers typically choose services that they expect to like (Hu et al. 2017). Reputation attribution appears when products receive unidimensional ratings describing multiple dimensions (e.g., value for money, appearance, durability). TripAdvisor acknowledges the issue of reputation attribution and offers reputation scores in four secondary dimensions (i.e., location, cleanliness, service, value). Even though such multidimensional systems better describe service quality, they do not generalize to an arbitrary set of dimensions, and they usually rely on a few dimensions that humans can efficiently rate. Finally, because some of the rated products or services are dynamic (e.g., venues evolve on TripAdvisor), their respective reputation systems likely experience reputation staticity.

Machine-based link analysis approaches require a network of interactions related to service quality in order to work (e.g., Jurczyk and Agichtein 2007). At the same time, they do not rely on any evaluation measure (either objective through testing or subjective through human raters). Hence, their applicability to contexts that require human-perceived service quality, such as online work, is limited. Applicability is also an issue for *hybrid-based* reputation systems (Section 2.1.3) that either focus on different objectives (e.g., node quality in peer-to-peer networks; see Damiani et al. 2002) or require information that is not freely available (e.g., social behavior; see Sabater and

Sierra 2001b) in online labor markets. Overall, these machine-based and hybrid approaches do not focus on addressing reputation inflation, reputation attribution, and reputation staticity, as these shortcomings do not pose significant limitations in their respective contexts (i.e., peer-to-peer networks and question-answering communities).

One could argue that given the contextual similarities between crowdsourcing settings and online labor markets, their reputation systems (Section 2.1.3) could be applicable in both contexts. Digital workplaces, however, differ from many crowdsourcing settings in that workers are highly paid and highly skilled (Paolacci et al. 2010, Kokkodis and Ipeirotis 2014). As a result, the objective of crowdsourcing systems to filter out adversarial workers does not apply to the focal context, as none of the high-skilled workers in online labor markets will purposely perform subpar work (e.g., by mindlessly labeling images, which is a typical task on AMT). Furthermore, low-skilled AMT workers are not susceptible to supply trends that often require reskilling; hence, skillset diversity and heterogeneity are not evident in crowdsourcing settings. As a result, crowdsourcing approaches do not focus on and do not solve reputation staticity,¹ reputation attribution, and reputation inflation.

Prior research has also proposed machine-based and hybrid reputation systems specifically for online labor markets (Christoforaki and Ipeirotis 2015, Daltayanni et al. 2015). Machine-based approaches rely on item response theory (Hambleton et al. 1991) to continuously generate test questions and evaluate the expertise of a user on a given skill (Christoforaki and Ipeirotis 2015). Many online labor markets already use such tests to certify workers on specific skills (Upwork 2018). In theory, these approaches could address reputation attribution, as they can evaluate the expertise of each worker on a given skill. In practice, however, they are costly (in terms of money and, most importantly, time), and they do not scale to multiple skills since they test one skill at a time (Upwork 2018). Even further, workers who take these tests have the option to either reveal or hide their scores, which results in the disclosure of positive-only certifications (Ipeirotis 2013). Finally, from a platform's perspective, creating and maintaining tests for hundreds of skills could incur additional costs.

A hybrid approach that is relevant to this work uses link-analysis and implicit reputation signals (e.g., "shortlisted," "hired," "ignored") to estimate a reputation score for each worker (WorkerRank; see Daltayanni et al. 2015). Specifically, by creating a historical graph between jobs and workers, the WorkerRank algorithm compares how different employers rank workers through their hiring processes. By construction, this approach does not focus and does not solve reputation

attribution and reputation staticity; however, because it has the potential to provide more representative reputation scores, it implicitly addresses reputation inflation. (Section 5 empirically shows the superiority of the proposed approach over the WorkerRank algorithm across multiple dimensions.)

2.2.3. Designing Dynamic Reputation Systems. Table 1 compares these relevant reputation systems found in academia and industry and identifies that they do not explicitly address reputation inflation, reputation attribution, and reputation staticity. This paper fills this gap by presenting the design principles that a reputation system needs in order to successfully attack reputation attribution, reputation staticity, and reputation inflation. Specifically, these principles are (1) decomposition of any arbitrary combination of skills into a set of finite competency dimensions, (2) dynamic estimation of competency-specific reputation, and (3) on-demand aggregation of these competency-specific reputations to estimate skillset-specific reputation scores. The first principle solves reputation attribution, as it allows for the efficient decomposition of reputation to a finite set of competency dimensions. The second principle allows for competency-specific dynamic estimation of quality and directly addresses reputation staticity. The third principle creates skillset-specific estimates of quality. Since human abilities follow normal-like distributions (Schmidt and Hunter 1983), accurate skillset-specific reputation scores address reputation inflation and facilitate worker differentiation.

2.3. Connection with Recommender Systems

Table 1 demonstrates the gap within the reputation systems literature that this study fills. Yet, and similar to reputation systems, recommender systems also resolve information asymmetries and allow customers to make better-informed decisions (Adomavicius and Tuzhilin 2005, Dellarocas 2006, Adamopoulos 2013). *Can such recommender systems resolve reputation staticity, reputation attribution, and reputation inflation and provide current skillset-specific worker-reputation scores in online labor markets?*

2.3.1. Differences Between Reputation and Recommender Systems. Reputation and recommender systems serve conceptually different objectives (Jøsang et al. 2013). Overall, reputation frameworks are broader in scope. For instance, reputation scores can generate product rankings (Amazon 2020), and, at the same time, they can provide quality estimates that affect product valuations (Moreno and Terwiesch 2014) and shape expectations (Ho et al. 2017). Even further, reputation systems can provide product managers with constructive feedback on how to improve their products (Proserpio and Zervas 2017). Besides, and

Table 1. Comparison of Relevant Literature on Reputation Systems

Paper or commercial implementation	Context	Applicable to OLMs	Reputation inflation	Reputation attribution	Reputation staticity	Requires testing	Methodology
Commercial reputation systems	OLMs	✓	×	×	×	×	Human-based, rating assignment
Multidimensional reputation systems	TripAdvisor	✓	×	+	×	×	Human-based, rating assignment
Christoforaki and Ipeirotis (2015)	OLMs	✓	⊕	⊕	⊕	✓	Machine-based, skill-specific testing (IRT)
Jurczyk and Agichtein (2007)	Q&A forums	×	×	×	×	×	Machine-based, network analysis
Zhang et al. (2007)	Q&A forums	×	×	×	×	×	Machine-based, network analysis
Bouguesia et al. (2008)	Q&A forums	×	×	×	×	×	Machine-based, network analysis
Kamvar et al. (2003)	Peer-to-peer sharing networks	×	×	×	×	×	Machine-based, network analysis
Sabater and Sierra (2001b)	E-commerce platforms, social settings	×	×	×	×	×	Hybrid, combines social, individual and ontological dimensions
Allahbakhsh et al. (2012)	Crowdsourcing (AMT)	×	×	×	+	×	Hybrid, graph-based analysis
Jagabathula et al. (2014)	Crowdsourcing (AMT)	×	×	×	×	×	Hybrid, graph-based analysis
Xie et al. (2015)	Crowdsourcing (AMT)	×	×	×	×	×	Hybrid, Bayesian game framework
Daltayanni et al. (2015)	OLMs	✓	?	×	×	×	Link analysis, implicit feedback
This research	OLMs	✓	✓	✓	✓	×	Deep learning, hidden Markov models

Notes. OLMs, online labor markets; IRT, item response theory; Q&A, question-answering communities; AMT, Amazon Mechanical Turk. The column “Applicable to OLMs” identifies whether the approach could be deployed in an online labor market. Columns “Reputation inflation,” “Reputation attribution,” and “Reputation staticity” identify whether the research addresses these shortcomings of current reputation systems in online labor markets. The column “Requires testing” identifies whether the approach requires workers (users) to take tests.

⊕: Testing has the potential to address reputation inflation, reputation attribution, and reputation staticity, but it requires investment in time (and money) and it has scalability and cost constraints (Section 2.2.2). +: Approach allows reputation in four dimensions; hence, to a certain degree, it addresses reputation attribution. Online Appendix E.3 shows that the proposed approach outperforms such multidimensional reputation systems. -: Approach incorporates the timing of each feedback in estimating reputation, so, in theory it addresses reputation staticity. ? : Approach could potentially address reputation inflation.

specifically to online labor markets, employers can use reputation scores to invite workers to apply to their job openings (Rahman 2018b). Once workers apply to a job opening, employers can rank and choose applicants according to their reputation scores (Kokkodis et al. 2015, Abhinav et al. 2017, Kokkodis 2018).

On the other hand, recommender systems often serve a single objective. For instance, they recommend products that customers usually buy together (basket analysis; see Agrawal et al. 1994). Or, they use customers' observed actions to recommend items that users will choose next (next-item recommendations; see Rendle et al. 2010, Quadrana et al. 2018). Or, they provide product recommendations based on what similar users with the targeted user have liked in the past (Billsus and Pazzani 1998, Breese et al. 1998, Delgado and Ishii 1999, Adomavicius and Tuzhilin 2005). Specifically to online labor markets, existing recommender systems (1) rank job applicants within a given opening (Kokkodis et al. 2015, Abhinav et al. 2017), (2) recommend tasks for workers to apply (Goswami et al. 2014, Baba et al. 2016, Horton 2017), and (3) provide career path recommendations (Patel et al. 2017, Kokkodis and Ipeirotis 2020).

In many cases, reputation and recommender systems work together to increase trust and reduce uncertainty (Jøsang et al. 2013). In online labor markets, in particular, reputation scores are often predictors in probabilistic recommender systems of different objectives. For instance, job-applicant recommenders use worker reputation as one of the attributes in their classification approaches (Kokkodis et al. 2015, Abhinav et al. 2017). Career-path recommenders use worker-reputation scores to provide relevant skill recommendations (Kokkodis and Ipeirotis 2020). Systems that recommend tasks to workers use reputation scores to identify the most appropriate assignments (Hossain and Arefin 2019).

2.3.2. Recommender Systems as Worker-Reputation Frameworks. Despite these conceptual differences, recommender systems can adjust to provide worker-reputation scores. Traditional recommenders predict the rating a user would assign to an item (Ricci et al. 2011, Adamopoulos and Tuzhilin 2014). To apply these systems in the focal context, we need to map ratings, users, and items to their respective entities in a worker-reputation framework for online labor markets. Since workers' reputation is the desired outcome of the framework, workers map to a recommender system's users.

The mapping of items and ratings is more complicated. In traditional recommender systems, the items are static entities that do not evolve over time (e.g., movies, smartphones, songs). Besides, multiple users of traditional recommender systems buy,

experience, and rate identical items (e.g., the same movie, song, smartphone). These multiple ratings per item are necessary for algorithms such as collaborative filtering to provide item recommendations that similar users to the focal user have liked in the past (Koren 2010, Adamopoulos and Tuzhilin 2013). In the focal context, the rated items are the completed tasks. In online labor markets, very rarely (if ever) two tasks are identical; even tasks that require the same skills might have different objectives. As a result, we do not observe multiple ratings for each task; instead, each task receives only a single rating. Hence, to transform recommender systems to a worker-reputation framework, we need to create static items for which multiple workers receive ratings.

A straightforward noisy transformation is to consider tasks that require the same skillsets as identical items. Then, the skillset-specific average feedback score that each worker receives can map to a recommended item's rating. Using these mappings, I can fill the user-item matrix and implement matrix-completion algorithms (Adomavicius and Tuzhilin 2005, Koren 2010) that will provide worker-reputation scores for each available skillset (item). Even though this mapping is not perfect (e.g., each worker might perform multiple tasks that require the same skillsets over time and receive different ratings), it provides skillset-specific reputation scores and potentially addresses reputation attribution. What about reputation staticity?

The rich literature on recommender systems provides a plethora of sequence-aware (or session-based) approaches (Hsueh et al. 2008, Zang et al. 2010, Jannach et al. 2015, Quadrana et al. 2018) that could potentially address reputation staticity. These systems explicitly model sequences of past events and recommend items according to the user's short-term behavior (Quadrana et al. 2018). Because they focus on next-item recommendations, sequence-aware approaches use implicit user feedback (i.e., whether a user has bought a product), ignoring explicit feedback ratings that describe whether the user actually liked the bought product (Devooght and Bersini 2017, Kula 2018, Quadrana et al. 2018). Yet, most reputation frameworks require explicit feedback ratings to work (Tadelis 2016, Luca 2017, Einav et al. 2016, Amazon 2018, eBay 2018, Stack Overflow 2018, Kokkodis and Lappas 2020, Kokkodis et al. 2020a). As a result, the application of next-item recommenders in the studied context will require significant encoding assumptions of user actions to generate next-item recommendations that provide skillset-specific worker-reputation scores. I discuss these encoding assumptions in detail in Section 5.4.1.

Table 2 compares relevant literature in recommender systems with the proposed approach and demonstrates the necessary modifications that relevant

Table 2. Comparison of Relevant Literature of Recommender Systems

Type of recommender systems	Objective	Required modifications	Reputation inflation	Reputation attribution	Reputation staticity	Methodology
Collaborative filtering (Adomavicius and Tuzhilin 2005, Koren et al. 2009)	Recommend items that similar users to the targeted user have liked in the past	Skillset \mapsto item worker \mapsto user avg. feedback \mapsto rating	✗	✗	✗	Singular value decomposition, slope one
Content-based (Adomavicius and Tuzhilin 2005)	Recommend items based on commonalities between items that a user has rated highly in the past	Skillset \mapsto item worker \mapsto user avg. feedback \mapsto rating	✗	✗	✗	Cosine similarity, Euclidean distance, predictive modeling
Basket analysis (Agrawal et al. 1994)	Recommend items that users often buy together	Not applicable	✗	✗	✗	A priori, association rule mining
Sequence-aware (next-item) recommenders (Quadrana et al. 2018)	Recommend an item that a user should buy next based on the observed user history	Encoding of user actions (implicit feedback) to represent worker reputation (Section 5.4.1)	✗	✗	✗	Convolutional neural networks, long short-term memory, factorized personalized Markov chain
Recommenders in OLMs (Abhinav et al. 2017, Kokkodis et al. 2015, Hossain and Arefin 2019)	Recommend skills and tasks to workers, or job applicants to employers	Predictive models adjusted to regression output	✗	✗	✗	Logistic regression, Bayesian networks, support vector machines, decision trees, random forest
This research	Provide current and skillset-specific worker reputation	None	✓	✓	✓	Deep learning, hidden Markov models

Notes. Avg. feedback, the average feedback score that a worker receives on a given skillset (item); OLMs, online labor markets. The column “Required modifications” summarizes the necessary modifications and assumptions that recommenders need to make in order to provide skillset-specific worker-reputation scores. The columns “Reputation inflation,” “Reputation attribution,” and “Reputation staticity” identify whether the recommender system adaptation addresses these shortcomings of current reputation systems in online labor markets.

✗: Required encoding and item assumptions could potentially address reputation attribution and reputation staticity. Section 5.4 empirically shows that these assumptions do not resolve reputation attribution and reputation staticity sufficiently, and, as a result, they also do not resolve reputation inflation.

recommender systems need to make to provide worker-reputation scores. When applied in Section 5.4 and Figure 8, these modifications significantly hurt the performance of various recommenders, highlighting the need for a context-appropriate reputation system that addresses reputation inflation, reputation attribution, and reputation staticity.

3. A Dynamic, Multidimensional Reputation Framework

Section 2.2.3 summarized the three design principles that a reputation framework should follow in order to provide current skillset-specific worker quality scores. Based on these principles, I structure a reputation framework (HMM-W2V framework) that consists of three components. Component A focuses on decomposing skills into competency dimensions that allow the framework to generalize and accommodate any arbitrary number of skills. Component B builds a dynamic model that combines multiple signals to estimate a worker’s current, competency-specific quality.

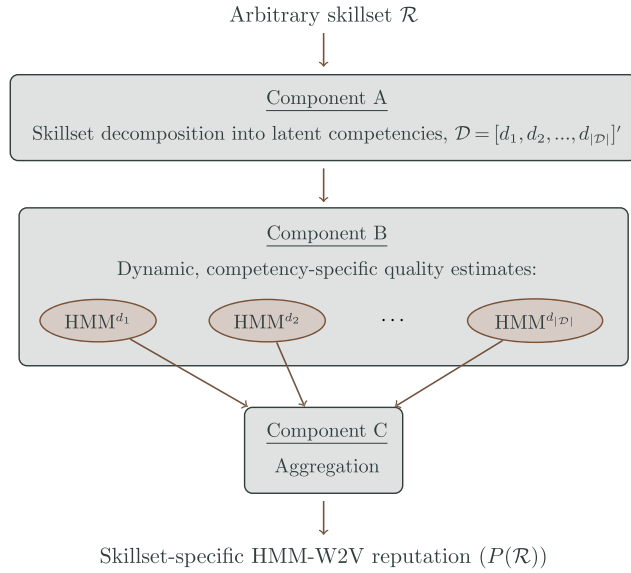
Component C aggregates competency-specific predictions to get a current estimate of the worker’s quality on any set of skills. Figure 1 draws the interconnections of these three components, which I describe in detail next.

3.1. Component A: Skills Decomposition

Workers can work on any arbitrary combination of skills, which creates a space of tens of thousands of available combinations of unique skillsets (Table 3). In theory, I could directly estimate the reputation of each worker on any observed skillset. This approach, however, has three drawbacks. First, independent of the size of the analyzed data, considering distinct skillset-specific observations will result in sparse training data sets. Second, such observations would ignore (by construction) correlations between various skillsets. Third, the entrance of new skills would require retraining of the framework.

To overcome these drawbacks, I use a distributed representation of words model (word embedding, W2V;

Figure 1. (Color online) The HMM-W2V Framework Provides Current, SkillSet-Specific Worker-Quality Estimates



Notes. The three design principles (Section 2.2.3) define the three-component structure of the HMM-W2V framework. Component A maps any set of skills to competency dimensions. Component B provides dynamic models (hidden Markov models, HMM) that make current, competency-specific quality estimates. Component C aggregates these estimates to provide a skillset-specific reputation.

see Mikolov et al. 2013) that projects individual skills into a set of competency dimensions. W2V embeds words from a vocabulary into a lower-dimensional space, in which semantically similar words appear close to each other, whereas semantically dissimilar words appear far away from each other (Mikolov et al. 2013). In the context of online work, a “skill” maps to a “word” and a “skillset” to a “document.” Based on this representation, W2V projects contextually similar skills close to each other in a $|\mathcal{D}|$ -dimensional space of competencies. (The actual number of competency dimensions $|\mathcal{D}|$ is a hyperparameter of the framework.)

The HMM-W2V framework maps any observed skillset \mathcal{R} into a $|\mathcal{D}|$ -dimensional vector space of

competencies (Figure 1). To do so, it averages competency-specific scores of each skill in a given skillset. In particular, a skillset \mathcal{R} maps to an aggregated W2V representation as follows:

$$\hat{w}_{\mathcal{R}}^d = \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} W2V^d(r), \quad \forall d \in \mathcal{D}, \quad (1)$$

where $\hat{w}_{\mathcal{R}}^d$ is the d -competency score of skillset \mathcal{R} , and $W2V^d(r)$ is the W2V score for skill r in dimension d . I normalize these weights to create a scale-invariant decomposition through a softmax transformation:

$$\begin{aligned} \mathbf{w}_{\mathcal{R}} &= \text{softmax}\left(\left[\hat{w}_{\mathcal{R}}^1, \hat{w}_{\mathcal{R}}^2, \dots, \hat{w}_{\mathcal{R}}^{|\mathcal{D}|}\right]'\right) \\ &= \left[w_{\mathcal{R}}^1, w_{\mathcal{R}}^2, \dots, w_{\mathcal{R}}^{|\mathcal{D}|}\right]'. \end{aligned} \quad (2)$$

Alternatively, a distributed memory model (D2V; see Le and Mikolov 2014) or simpler clustering approaches could map skillsets into vectors of real numbers. Furthermore, in addition to the skillsets, W2V or D2V could also consider the job-description text. Online Appendix C.1 and Figure 11(a) discuss and empirically compare these alternative approaches.

3.2. Component B: Competency-Specific Dynamic Quality Assessment

At any given time, each worker has a quality level in each competency dimension $d \in \mathcal{D}$. This quality is latent (unobserved). However, for each completed task t with required skills \mathcal{R} , the market observes a feedback score (Y_t) that the worker receives. This feedback maps into competency-specific scores through the weighting vector $\mathbf{w}_{\mathcal{R}}$: $Y_t^d = w_{\mathcal{R}}^d Y_t$, $\forall d \in \mathcal{D}$. These scores (Y_t^d) form a sequence of proxies of the worker's quality in each competency dimension $d \in \mathcal{D}$.

In addition to latent, a worker's quality dynamically evolves. Specifically, a worker's quality can change in any competency dimension, either by gaining experience on the platform and better understanding the expectations of the employers or by learning new skills and continuously expanding current knowledge and abilities. The HMM-W2V framework formulates

Table 3. Data Overview

	Observations	Mean	Median	Standard deviation	Min	Max
Skills per worker	662,423	9.4	9	5.8	1	61
Tasks per worker	58,459	6.3	5	3.6	1	59
Skills per task	58,459	3.2	3	2.5	1	38
Task compensation (\$)	58,459	170	31	1,042	3	70,218
Task applications	662,423					
Completed tasks	58,459					
Unique skills (\mathbb{R})	547					
Unique observed skillsets	17,563					

Notes. Workers work remotely from 141 countries. Data span 12 months.

this evolution through a hidden Markov model (HMM): each worker operates from a latent, competency-specific state, which determines the worker's propensity to perform with score Y_t^d . Every time a worker completes a new task and receives a new feedback score, the framework observes new evidence about the worker's competency-specific qualities and stochastically transitions the worker to new latent states. The framework assumes a set of \mathcal{S}^d latent states that describe K^d different levels of quality for each competency d , $\mathcal{S}^d = \{s_1, s_2, \dots, s_{K^d}\}$.

HMM Structure. Every new worker who joins the platform has an unknown quality across the competency dimensions $d \in \mathcal{D}$. As the worker completes new tasks on the platform, the market observes signals (i.e., through task outcomes) that correlate with the latent worker quality. To capture this behavior, the HMM-W2V framework assumes an initial latent state s_1 , where all new workers land. This state makes an average initial estimate of workers' competency-specific quality.² Once the workers complete their first task and emit an observation (i.e., $Y_1^d, \forall d \in \mathcal{D}$), they stochastically transition to different states according to the parameters of the model.

To define an HMM for a given competency dimension d , I need (1) a vector of initial state probabilities π^d , (2) a transition matrix T^d that stores the transition probabilities between states, and (3) an emission matrix E^d that describes the state-specific probability distributions for observations Y_t^d . Since every new worker lands in state s_1 , the initial probability vector of each HMM is the following:

$$\pi^d = [1, 0, 0, \dots, 0]'. \quad (3)$$

A worker's history provides multiple observable signals that correlate with transitions to new quality states (e.g., total wages received, hiring rates, number of completed tasks). Such historical attributes define a vector \mathbf{Z}_{t-1} .³ By weighing this vector with $w_{\mathcal{R}}^d$, the framework forms vectors $\mathbf{Z}_{t-1}^d = w_{\mathcal{R}}^d \mathbf{Z}_{t-1}$ that capture competency-specific histories. Each \mathbf{Z}_{t-1}^d directly affects the transition probabilities to different states (i.e., matrix T^d). Formally, assume that a given worker completes task $t - 1$ from state s_k in a given dimension d . Once the framework observes the outcome of task $t - 1$, it estimates the transition probability of this worker to move to state s_l as follows:

$$\begin{aligned} \lambda_{\gamma_{kl}^d \mathbf{Z}_{t-1}^d}^{d, s_k s_l} &= \Pr\left(S_t^d = s_l | S_{t-1}^d = s_k; \gamma_{kl}^d, \mathbf{Z}_{t-1}^d\right) \\ &= g^d(\gamma_{kl}^d \mathbf{Z}_{t-1}^d). \end{aligned} \quad (4)$$

In the previous equation, γ_{kl}^d is the vector of coefficients of state s_k that define the weights of \mathbf{Z}_{t-1}^d in estimating the transition probability to state s_l .

Function g^d transforms the product $\gamma_{kl}^d \mathbf{Z}_{t-1}^d$ into a probability. (The choice of function g^d is context- and data-specific. I discuss this in Section 5.1 and Online Appendix A.) The complete transition matrix for a given worker after completing task $t - 1$ in dimension d is as follows:

$$T^d(\Gamma^d, \mathbf{Z}_{t-1}^d) = \begin{bmatrix} \lambda_{\gamma_{11}^d \mathbf{Z}_{t-1}^d}^{d, s_1 s_1} & \lambda_{\gamma_{12}^d \mathbf{Z}_{t-1}^d}^{d, s_1 s_2} & \dots & \lambda_{\gamma_{1K^d}^d \mathbf{Z}_{t-1}^d}^{d, s_1 s_{K^d}} \\ \lambda_{\gamma_{21}^d \mathbf{Z}_{t-1}^d}^{d, s_2 s_1} & \lambda_{\gamma_{22}^d \mathbf{Z}_{t-1}^d}^{d, s_2 s_2} & \dots & \lambda_{\gamma_{2K^d}^d \mathbf{Z}_{t-1}^d}^{d, s_2 s_{K^d}} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_{\gamma_{K^d 1}^d \mathbf{Z}_{t-1}^d}^{d, s_{K^d} s_1} & \lambda_{\gamma_{K^d 2}^d \mathbf{Z}_{t-1}^d}^{d, s_{K^d} s_2} & \dots & \lambda_{\gamma_{K^d K^d}^d \mathbf{Z}_{t-1}^d}^{d, s_{K^d} s_{K^d}} \end{bmatrix}, \quad (5)$$

where $\Gamma^d = [\gamma_{11}^d, \gamma_{12}^d, \dots, \gamma_{K^d K^d}^d]'$.

Similarly, observed worker characteristics (e.g., hourly rate, average feedback score) are correlated with the observed emissions of the HMM (matrix E^d). These characteristics form a vector \mathbf{X}_t . The framework makes this vector competency-specific by weighting its elements with $w_{\mathcal{R}}^d$; that is, $\mathbf{X}_t^d = w_{\mathcal{R}}^d \mathbf{X}_t$. Formally, the conditional probability of observing Y_t^d given the current state of the worker S_t^d is

$$\Pr(Y_t^d | S_t^d = s_k; \theta_k^d, \mathbf{X}_t^d) = f^d(\theta_k^d \mathbf{X}_t^d), \quad (6)$$

where $f^d(\cdot)$ is a continuous probability distribution (e.g., Beta), and θ_k^d is the parameter vector of the continuous distribution for state k and competency d . The complete parameter vector Θ^d (for all states $s_k \in \mathcal{S}^d$) is as follows:

$$\Theta^d = [\theta_1^d, \theta_2^d, \dots, \theta_{K^d}^d]'. \quad (7)$$

Online Appendix B presents the derivation of the likelihood and the subsequent process of estimating the parameters of the model. Online Appendix C discusses the choice of emission (g^d) and transition (f^d) functions, as well as the tuning of the total number of states K^d . Finally, Online Appendix C compares alternative approaches for modeling component B of the HMM-W2V framework, including recurrent neural networks and gradient boosting.

3.3. Component C: Aggregation

This process happens independently for each competency dimension $d \in \mathcal{D}$. As a result, for a given worker i who has completed t tasks, each HMM estimates a current competency-specific quality p_{it}^d [i.e., a stochastic draw from the continuous emission distribution of Equation (6)]. To estimate the quality of worker i for any given set of skills \mathcal{R} , the HMM-W2V framework aggregates the available competency-specific estimates:

$$P_{it}(\mathcal{R}) = \sum_{d \in \mathcal{D}} p_{it}^d. \quad (8)$$

Summation of these estimates allows each competency to contribute to the skillset-specific reputation by its respective, skillset-specific weight [recall that Equation (2) softmaxes the weights of each dimension]. Online Appendix C and Figure 11(c) compare alternative aggregating approaches.

4. Data Description and Model-Free Evidence

I build and evaluate a version of the proposed framework on a set of real transactions from a major online labor market, LaborBazaar (pseudonym). The focal data form a snapshot of 662,423 task applications that led to 58,459 completed tasks by 13,510 workers. LaborBazaar supports diverse tasks from different categories, including software and web development, writing, sales, marketing, and data science. Table 3 summarizes the data set. Overall, 547 unique skills create a total of 17,563 unique skillsets. Workers participate in this platform remotely from 141

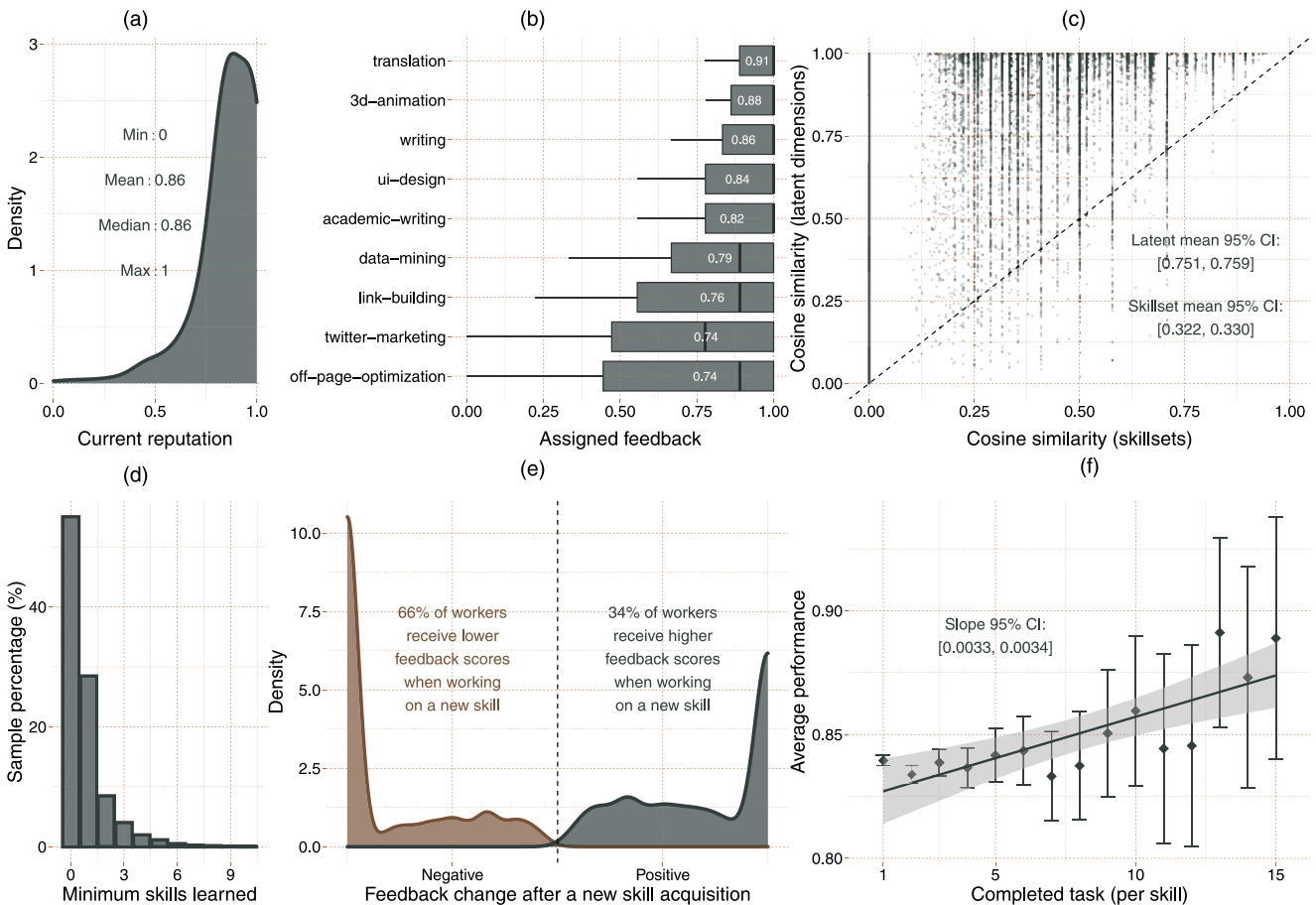
different countries. I follow the actions of these workers for 12 consecutive months.

4.1. Model-Free Evidence

On LaborBazaar, employers rate workers after the completion of a task with a score $Y \in \{0, 1/9, 2/9, \dots, 9/9\}$. The platform supports a state-of-the-art reputation system (Filippas et al. 2018): The employer and the worker each get two weeks to leave their feedback score. This is a double-blind process, where neither party learns its rating before leaving a rating for the other party. Figure 2(a) shows the resulting feedback distribution of the focal data. The mean and median of this distribution are both 0.86. Most of the workers appear to perform almost perfectly (reputation inflation).

At the same time, the heterogeneity in terms of skills and qualifications on the platform in combination with the fact that feedback scores are assigned uniformly (reputation attribution) adds noise to the

Figure 2. (Color online) Model-Free Evidence of Shortcomings of Current Reputation Systems.



Notes. The six panels describe the reputation system of LaborBazaar. Panel (a) shows that the accumulated feedback scores of workers are inflated. Panel (b) shows that different skills have different feedback score distributions. Panel (c) shows that workers work on heterogeneous tasks—skillset average cosine similarity between consecutive tasks is 0.33—highlighting reputation attribution. Panel (d) shows that workers evolve by learning new skills. Such new skill acquisitions usually initially hurt worker reputation [panel (e)]. However, over time, workers gain experience and perform better [panel (f)]. As a result, panels (d), (e), and (f) combined highlight the problem of reputation staticity.

inflated distribution. Figure 2(b) shows that skills accumulate different feedback scores: from translation, with median 1 and mean 0.91, to Twitter marketing with median 7/9 and mean 0.74. Figure 2(c) further shows that workers work on heterogeneous tasks. Specifically, the x -axis shows the cosine similarity of consecutive tasks in terms of skillsets, and the y -axis shows the cosine similarity between consecutive tasks in terms of projected latent dimensions (Section 3.1). The mean cosine similarity of consecutive tasks is 0.33, suggesting that workers indeed work on heterogeneous tasks (reputation attribution). The same figure further shows how mapping skillsets through W2V captures contextual similarities between different skillsets and, as a result, yields higher cosine similarity between consecutive skillsets (average mean similarity of consecutive tasks in the latent space is 0.75).

Figure 2(d) shows that around 47% of the workers of this sample use in the marketplace at least one new skill. Recall that I follow the focal workers only for 12 months. As a result, the observed skillset evolution happens during these 12 months. Figure 2(e) shows that when workers acquire a new skill, they often (66%) initially receive lower feedback scores. However, as they gain experience, Figure 2(f) shows that their reputation increases. These three graphs highlight the dynamic nature of workers, which suggests that reputation systems should adjust for reputation staticity.

4.2. Emission and Transition Variables

The HMM-W2V framework requires a set of emission and transition variables that describe vectors \mathbf{Z}_{t-1} and \mathbf{X}_t . Recall that each latent state represents a different distribution of expected service quality. Transitions between states are subject to the accumulated experience and feedback scores of each worker that form vector \mathbf{Z}_{t-1} . In particular, I allow five signals to affect transitions: (1) the current accumulated reputation of the worker, (2) the total money earned on the platform, (3) the total number of completed jobs, (4) the total

number of hours worked, and (5) the worker's hiring rate.

Similarly, vector \mathbf{X}_t that captures observed worker characteristics affects the emission probabilities. Such characteristics include the current reputation of the worker and the worker's current hourly rate. Table 4 shows the descriptive statistics of the variables that form vectors $\mathbf{X}_t, \mathbf{Z}_{t-1}$, as well as the outcome variable Y_t . I log-transform variables with long tails and standardize all nonbinary variables for faster convergence. Finally, note that this illustrative list of variables that formulate HMM transitions and emissions is context-specific. Online Appendix E.3 shows how alternative contexts require different variable choices.

5. Evaluation of the HMM-W2V Framework

The next paragraphs describe the modeling choices and hyperparameter tuning of the HMM-W2V framework and compare its performance with various alternative advanced reputation systems and modified recommender systems. I split the data into 10 folds that consist of different workers (i.e., each worker's complete history appears only in 1 of the 10 folds). I use 10-fold cross-validation to evaluate and compare each alternative design approach.

5.1. Design Choices and Hyperparameter Tuning

Alternative design options could model each component of the framework (Section 3). To identify the best design choices for the context of this study, I follow a grid-search approach. Specifically, I compare alternative modeling choices for components A, B, and C, and test the framework's performance under various numbers of dimensions $|D|$. Furthermore, I evaluate various combinations of numbers of states K^d , choices of transition functions g^d , and choices of emission functions f^d . Online Appendix C discusses the details of this grid-search approach.

Based on this analysis, the combinations that performed best in the focal context are the following:

Table 4. Descriptive Statistics for the Attributes in Vectors $\mathbf{X}_t, \mathbf{Z}_{t-1}$, and the Outcome Variable Y_t

	Mean	Median	Standard deviation	Min	Max
Observed outcome (Y_t)	0.85	1	0.25	0	1
Transition vector \mathbf{Z}_{t-1}					
Current reputation	0.86	0.86	0.15	0	1
Total money earned (\$)	4,289	459	11,735	0	152,526
Completed jobs	9.71	2	20.68	0	401
Work-hours	468	31	1,368	0	36,457
Hiring rate	0.07	0.05	0.07	0	1
Emission vector \mathbf{X}_t					
Current reputation	0.86	0.87	0.15	0	1
Hourly rate (\$)	11.23	8.69	13.21	3	397

◊ *Modeling component A*: W2V performs better than D2V and Gaussian mixture models [Online Appendix C, Figure 11(a)]. Furthermore, including the job-description text in W2V adds noise and does not improve the performance of component A [Online Appendix C, Figure 11(a)]. Hence, for the focal data set, I choose W2V.

◊ *Modeling component B*: Modeling dynamic transitions through a hidden Markov model performs significantly better ($p < 0.001$) than linear models, support vector machines, recurrent neural networks, and gradient boosting [Online Appendix C, Figure 11(b)]. Hence, I choose the HMM structure described in Section 3.2 for the focal data set.

◊ *Modeling component C*: Aggregating dimension-specific feedback according to Equation (8) performs on par or better than alternative aggregation approaches [Online Appendix C, Figure 11(c)]. Hence, I use Equation (8) for the rest of the analysis.

◊ *Number of dimensions*: $|\mathcal{D}| = 10$ performs better than alternative values [Online Appendix C, Figure 11(d)].

◊ *HMM parameters*: For the focal context I choose (Online Appendix C.5):

- ◊ $K = [3, 4, 4, 4, 4, 4, 3, 4, 3, 3]$.
- ◊ $f^d = \text{Beta } \forall d \in \mathcal{D}$.
- ◊ $g^d = \text{Multinomial logit } \forall d \in \mathcal{D}$.

Using these design choices and hyperparameter values, the HMM-W2V framework estimates the reputation of each available worker on any arbitrary set of skills.

5.2. Alternative Reputation Systems

Alternative approaches could also generate reputation scores for each available worker. Recent advances in machine learning provide multiple approaches that could model sequential observations and capture dynamic behavior. Furthermore, applicable context-specific approaches (e.g., WorkerRank; see Daltayanni et al. 2015) could potentially address reputation inflation (Table 1). To benchmark the performance of the HMM-W2V framework against such advanced alternative models, I implement and compare the following reputation systems:

◊ *Current reputation*: The accumulated feedback score of LaborBazaar. This score is a result of a human-based reputation system, where each worker gets rated upon completion of a task, and these ratings accumulate to form worker reputation.

◊ *Machine learning approaches*: Alternative hybrid reputation systems that combine (1) human input (ratings), (2) observed characteristics, and (3) alternative modeling choices. These alternative systems model the relationship

$$Y_t \sim G(\mathbf{Z}_{t-1}, \mathbf{X}_t), \quad (9)$$

where G represents the following:

◊ *Linear model*: G linearly regresses the dependent variable on $\mathbf{Z}_{t-1}, \mathbf{X}_t$.

◊ *Recurrent neural networks (LSTM)*: G captures the relationships between vectors $\mathbf{Z}_{t-1}, \mathbf{X}_t$, and Y_t through long short-term memory networks (LSTM; see Hochreiter and Schmidhuber 1997).

◊ *Gradient boosting regression (XGBoost)*: G captures the relationships between vectors $\mathbf{Z}_{t-1}, \mathbf{X}_t$, and Y_t through gradient boosting regression (Chen and Guestrin 2016).

◊ *Support vector regression (SVM-reg)*: G captures the relationships between vectors $\mathbf{Z}_{t-1}, \mathbf{X}_t$, and Y_t through an SVM regression model (Smola and Schölkopf 2004).

◊ *WorkerRank*: An advanced reputation system for online labor markets that uses implicit feedback from employers (e.g., “hired,” “invited”) to rank workers through a link analysis approach (Daltayanni et al. 2015).

5.3. Results

The next paragraphs benchmark the performance of the HMM-W2V framework against alternative reputation systems in terms of (1) ranking workers, (2) generating a representative reputation distribution, (3) estimating the service quality of the most dynamic “nonperfect” workers, and (4) ranking applicants within openings.

Ranking Workers. The ultimate goal of any reputation system is to generate *rankings* of products or services (in this case, workers) according to their expected service quality. Hence, accurate assessment of quality should rank workers according to their likelihood of performing well on any given skillset \mathcal{R} . I capture the ranking performance of each reputation system through the following measures:

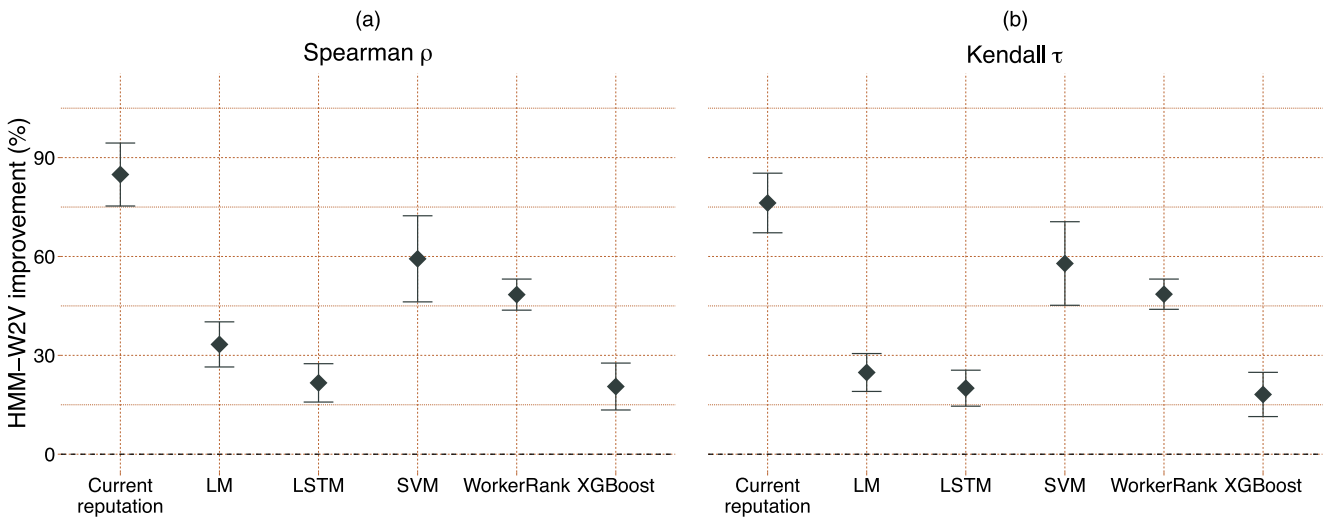
◊ *Ranking correlations*: Two ranking correlation coefficients (Kendall τ and Spearman ρ ; see Kendall 1938, Spearman 1904) test the ordinal associations between quality estimates and observed outcomes.

◊ *Ranking performance*: Detailed analysis of the average performance of workers ranked in different cohorts tests whether workers ranked in the top tiers perform consistently better than workers ranked in the bottom tiers.

◊ *Average lift*: Lift analysis estimates how much better top-ranked workers perform than bottom-ranked workers.

Figure 3 compares the proposed approach with alternative reputation systems in terms of ranking correlations. The y -axis shows the 10-fold cross-validated percentage improvement of the HMM-W2V framework over the x -axis reputation systems,

Figure 3. (Color online) Ranking Correlations of the HMM-W2V Framework and Alternative Reputation Systems



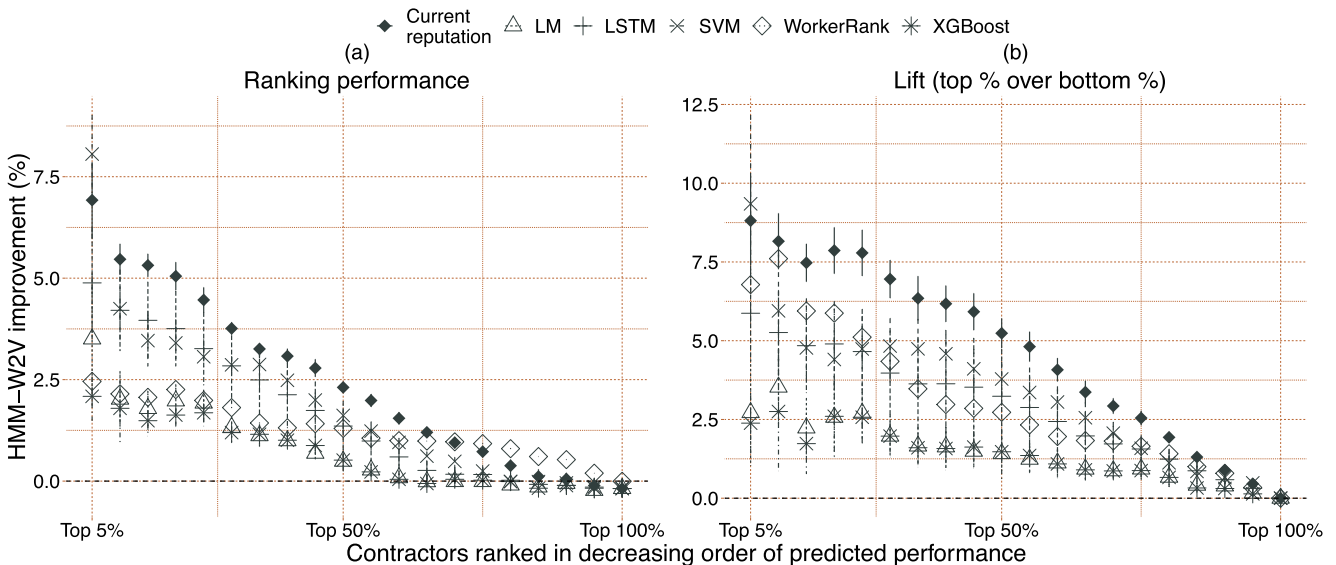
Notes. The y -axis shows the percentage improvement of the HMM-W2V framework over the x -axis reputation systems, in terms of Spearman ρ and Kendall τ . The HMM-W2V framework significantly ($p < 0.001$) outperforms all alternative reputation systems, with average 10-fold cross-validation improvements ranging between 20% and 85%. Error bars represent 95% confidence intervals.

in terms of Spearman ρ and Kendall τ . These ranking correlations capture how aligned each reputation system's ranking is with the observed performance of the workers. Across both metrics, the proposed approach significantly outperforms both current and alternative reputation systems: The HMM-W2V framework yields, on average, ~85% better rankings than the current reputation scores. At the same time, it significantly ($p < 0.001$) yields better rankings than all

alternative reputation systems (average improvement between 20% and 60%). Hence, the ordering of workers according to their predicted HMM-W2V reputation is significantly more accurate than their orderings according to alternative reputation systems.

Figure 4, (a) and (b), shows the ranking performance and lift improvement of the proposed approach over alternative reputation systems. The x -axis ranks workers according to their predicted reputation

Figure 4. (Color online) Ranking Performance and Lift of the HMM-W2V Framework and Alternative Reputation Systems



Notes. In both panels, the x -axis ranks workers according to their estimated reputation. The y -axis shows the 10-fold cross-validated percentage improvement of the HMM-W2V framework over each alternative reputation system. Panel (a) shows the ranking performance: Top-ranked workers according to the HMM-W2V reputation clearly outperform ($p < 0.001$) top-ranked workers from all alternative reputation systems by up to 8%. Panel (b) shows the average lift of each ranked cohort. The HMM-W2V framework yields up to 9% ($p < 0.001$) higher average lifts than the three baselines. Error bars represent 95% confidence intervals.

scores by each alternative approach. The y -axis shows the HMM-W2V reputation improvement over each alternative approach (a) in terms of the observed performance of the top-ranked workers, and (b) in terms of lift (i.e., how much better top-ranked workers perform than bottom-ranked workers). In Figure 4(a), the HMM-W2V framework significantly outperforms (up to 8%, $p < 0.001$) all alternative reputation systems. According to HMM-W2V reputation, higher-ranked workers yield significantly higher average performance. As the ranking moves from top to bottom, the improvement of the HMM-W2V reputation converges to zero as the top-ranked cohort includes a larger portion of the available contractors (i.e., the top-ranked sample's average performance converges to the population's average performance). Note that there is not a single point on the x -axis where an alternative reputation system outperforms the HMM-W2V framework. Similarly, in Figure 4(b), the proposed framework yields a higher average lift (up to 9%, $p < 0.001$) than all alternative approaches. This means that the rates of the HMM-W2V reputation top-ranked worker performance over the bottom-ranked worker performance are significantly ($p < 0.001$) higher than the respective rates of all alternative reputation systems.

Reputation Distribution. Figure 5 compares the resulting reputation distributions of the HMM-W2V framework and all alternative reputation systems. The focus is on estimating how close each resulting distribution is to the normal distribution. A normal distribution is more likely to represent the skillset-specific abilities of workers (Schmidt and Hunter 1983) and facilitate worker differentiation. The total

variation distance (Huber 2011, Kohl 2019) between any two distributions captures how close these distributions are. Hence, I use this distance to measure the closeness of each alternative reputation distribution to the normal distribution. Figure 5 shows that the proposed approach yields reputation distributions that are up to 37% closer to a normal distribution ($p < 0.001$) compared with the resulting distributions of alternative reputation systems.

Performance Evaluation on “Nonperfect” Workers. Many workers in these markets always appear to perform well, in part due to reputation inflation (Table 4). As a result, models can have high accuracy by always correctly predicting these “perfect” workers. The real challenge is for algorithms to accurately predict the performance of workers who occasionally underperform (i.e., minority class prediction performance; see Longadge and Dongre 2013). This is important for the market, as early identification of potential underperforming workers could prevent employers from having a disappointing experience (Section 6).

To test the performance of each alternative reputation system on “nonperfect” workers, I estimate ranking correlations (Spearman ρ , Kendall τ) for the subset of workers who receive at least one imperfect (< 1) feedback score. Figure 6 shows the results. The focal framework significantly ($p < 0.001$) outperforms all alternative reputation systems in the populations that are harder to predict and are often costly for the marketplace (Section 6).

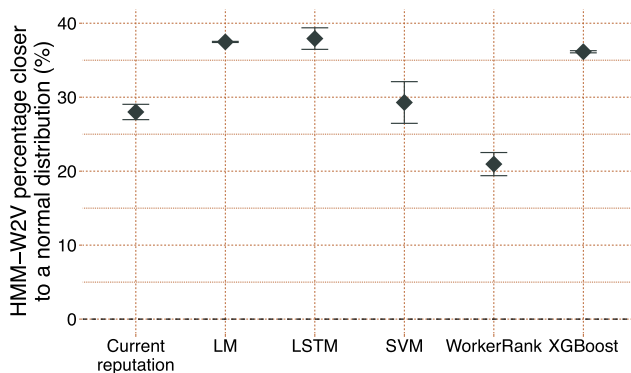
Performance Evaluation on Within-Opening Rankings.

The previous paragraphs demonstrate the superiority of the proposed approach compared with alternative reputation systems in terms of differentiating workers (ranking correlations, reputation distribution). Each reputation system, however, could generate within-opening rankings of *applicants* (i.e., workers who have applied for the job). Such within-choice-set rankings often affect buyer decisions (Ghose et al. 2012, 2014; Kokkodis et al. 2015). To compare the set of alternative reputation systems in terms of within-opening rankings, I do the following:

- For each reputation system, rank applicants within openings according to their reputation score.
- For each $n \in \{1, 2, 3, 4, 5\}$, observe the average performance of workers who got hired while ranking at the Top- n according to each reputation system.
- In the end, compare the average performance of each reputation system at Top- n .

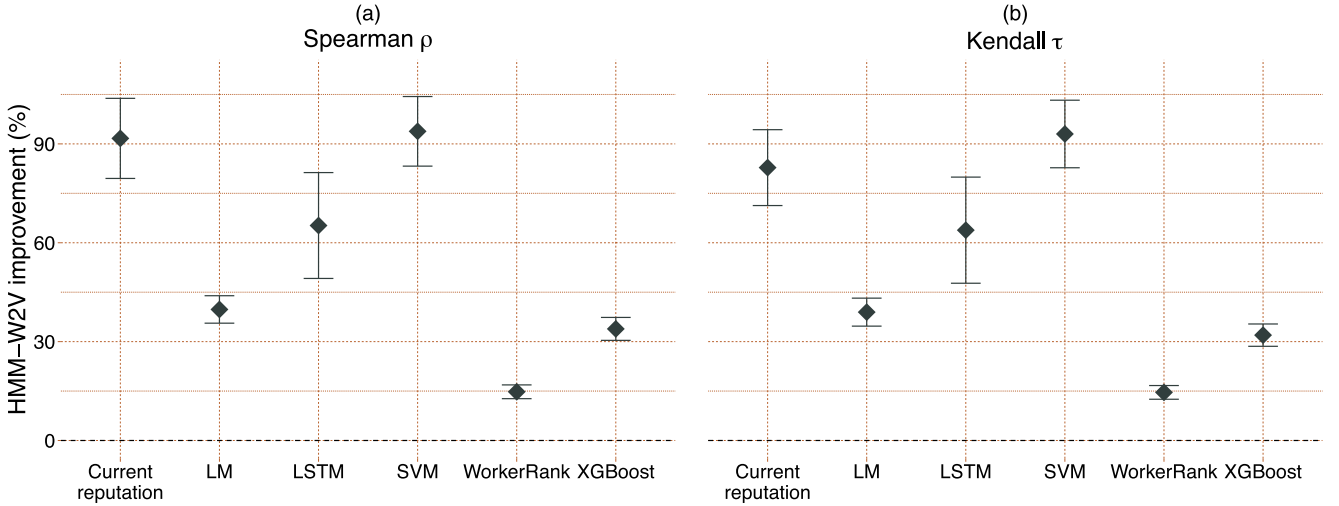
Figure 7 shows the results. The y -axis indicates the average performance of workers that each algorithm ranks at Top- n . The x -axis lists all reputation systems. The HMM-W2V framework outperforms ($p < 0.05$) all reputation systems but WorkerRank and XGBoost

Figure 5. (Color online) Reputation Distributions of the HMM-W2V Framework and Alternative Reputation Systems



Notes. The y -axis shows the 10-fold cross-validated improvement of the HMM-W2V framework in terms of the total variation distance compared with a normal distribution. Compared with alternative reputation systems, HMM-W2V reputation is up to 37% closer to a normal distribution. Error bars represent 95% confidence intervals.

Figure 6. (Color online) Evaluation on “Nonperfect” Workers



Notes. The y -axis indicates the improvement of the HMM-W2V framework over the x -axis reputation systems in terms of Spearman ρ and Kendall τ . The HMM-W2V framework outperforms all alternative reputation systems ($p < 0.001$) in identifying “nonperfect” workers. Error bars represent 95% confidence intervals.

across all n . It further outperforms WorkerRank at $p < 0.05$ for $n \in \{4, 5\}$, and at $p < 0.1$ for $n \in \{1, 2, 3\}$, and it outperforms XGBoost at $p < 0.1$ for $n = 1$ and at $p < 0.05$ for $n \in \{2, 3, 4, 5\}$. These results suggest that the HMM-W2V framework can offer better within-opening rankings of applicants than alternative reputation systems.

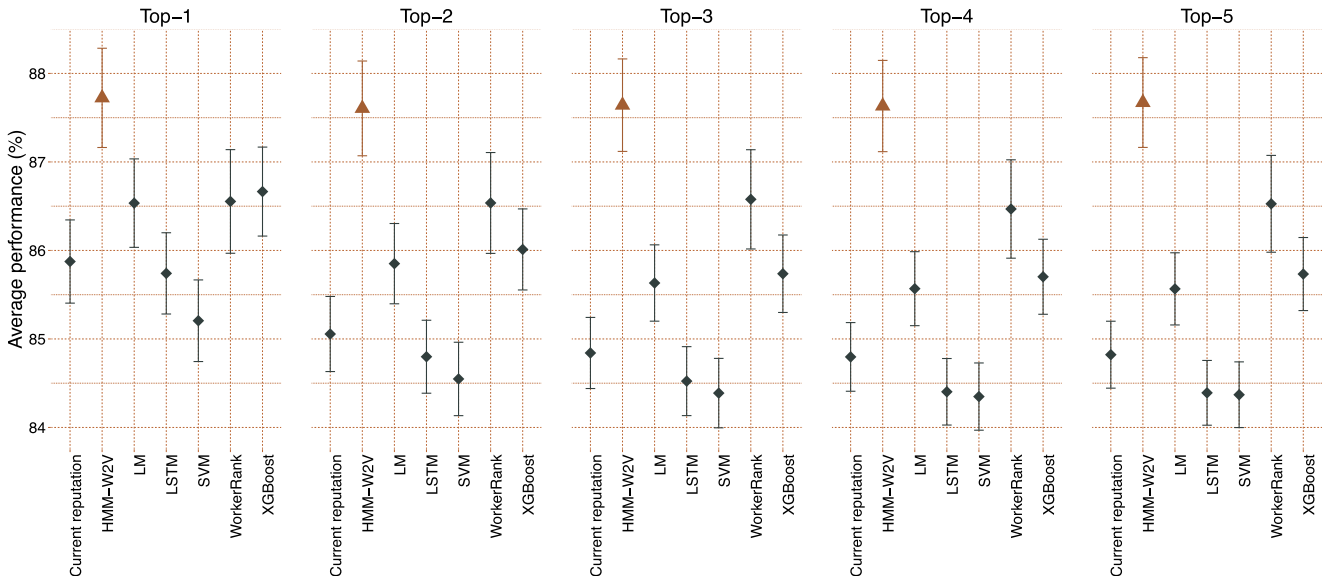
5.4. Empirical Evaluation of Recommender Systems

Section 2.3 summarized the conceptual differences between reputation and recommender systems. Table 2 further identified the necessary modifications and assumptions that recommender systems need to

make in order to provide worker-reputation scores. Section 5.4.1 empirically tests such adaptations of recommender systems and shows that they significantly underperform compared with the HMM-W2V framework. Finally, Section 5.4.2 illustrates how reputation and recommender systems can work together to enhance the transaction efficacy of a marketplace.

5.4.1. Adaptations of Recommender Systems as Reputation Frameworks. Section 2.3.2 and Table 2 identified that in order to adjust recommender systems to

Figure 7. (Color online) Comparison of Alternative Reputation Systems on Within-Opening Rankings



Notes. The y -axis indicates the average performance of the Top- n hired workers according to each reputation system. The HMM-W2V framework outperforms ($p < 0.05$) all reputation systems but WorkerRank and XGBoost across all n . It further outperforms WorkerRank at $p < 0.05$ for $n \in \{4, 5\}$, and at $p < 0.1$ for $n \in \{1, 2, 3\}$, and XGBoost at $p < 0.1$ for $n = 1$, and at $p < 0.05$ for $n \in \{2, 3, 4, 5\}$. Error bars represent 95% confidence intervals.

provide worker-reputation scores, I need to make the following assumptions:

- Skillset \mapsto item.
- Worker \mapsto user.
- Average worker's skillset-specific feedback score \mapsto rating.

Based on these, I transform and implement the following popular recommender systems:

- *Collaborative filtering*: Collaborative filtering recommenders are among the most powerful and widespread industry approaches (Adomavicius and Tuzhilin 2005, Su and Khoshgoftaar 2009, Meyer 2012, Adamopoulos and Tuzhilin 2015). For the focal context, I customize three popular powerful collaborative filtering frameworks: k -nearest neighbors (NN, Kantor et al. 2011), singular value decomposition (SVD, Kantor et al. 2011), and slope one (Lemire and Maclachlan 2005).

- *Neural network sequence recommender systems*: The rise and popularity of neural networks have motivated approaches that use such networks to build deep recommender systems. Given the nature of the focal context and the fact that workers evolve dynamically, a sequential recommendation model could capture latent and dynamic relationships by treating recommendations as a sequential prediction problem (Kung-Hsiang 2018). Convolutional neural networks (CNN) often model such sequential recommender systems that provide next-item recommendations (Kula 2018, Quadrana et al. 2018). As mentioned in Section 2.3.2, because these approaches focus on implicit feedback (i.e., whether a user chooses an item next), I need to encode the observed sequences of worker skillset-specific ratings into user actions. To do so, I combine a task's required skillset along with its respective observed performance. For instance, if a

worker completes a sequence of tasks that require {java} and {python} and receives feedback scores 8/9, and 1, then the sequence will be [{java}-8/9, {python}-1]. These transformations generate sequences of observations that a deep recommender system can use to predict the next skillset-score combination of each worker. I use these next-skillset-score predictions to infer a worker's skillset-specific reputation.

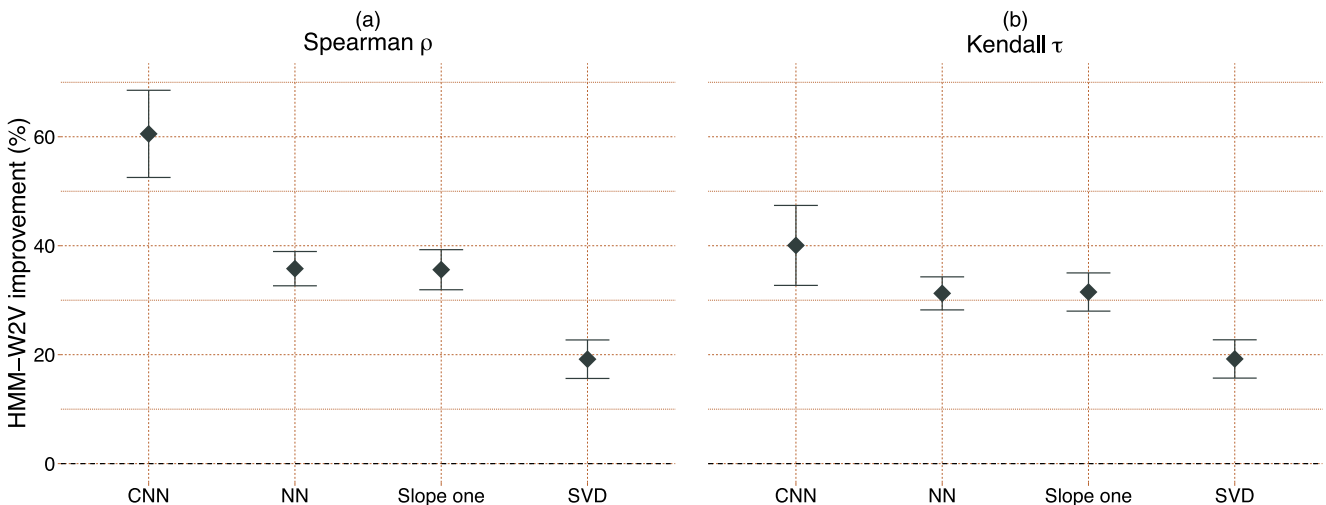
Through a grid search approach, I tune these models, and compare their performance with the HMM-W2V framework. Figure 8 shows how the proposed approach significantly ($p < 0.001$) outperforms all recommender systems in terms of ranking correlations. This underperformance shows that the necessary assumptions that transform recommender systems into worker-reputation frameworks likely hurt their performance. As a result, recommender systems do not sufficiently address reputation staticity, reputation attribution, and reputation inflation in online labor markets.

5.4.2. Collaboration of Reputation and Recommender Systems.

Reputation and recommender systems can work together to increase transaction efficacy. Specifically, worker-reputation scores can enhance the performance of recommender systems in online labor markets (Section 2.3.1). To demonstrate, I build existing job-applicant recommenders (Kokkodis et al. 2015, Abhinav et al. 2017), and I test their performance with and without using HMM-W2V reputation. Specifically, I build models that estimate the hiring probability of each job applicant (Kokkodis et al. 2015, Abhinav et al. 2017):

$$\Pr(\text{Hire}|\mathbf{W}) = h(\mathbf{W}), \quad (10)$$

Figure 8. (Color online) The HMM-W2V Framework Outperforms Adaptations of Recommender Systems



Notes. The y-axis shows the percentage improvement of the HMM-W2V framework over the x-axis recommender systems, in terms of Spearman ρ and Kendall τ . The HMM-W2V framework significantly ($p < 0.001$) outperforms all alternative recommender systems, with average 10-fold cross-validation improvements ranging between 20% and 60%. Error bars represent 95% confidence intervals.

where $h \in \{\text{Logistic regression, Bayesian networks, Random forest, Gradient boosting, Neural networks}\}$. The vector of observed job-applicant characteristics \mathbf{W} takes the following forms:

$$\mathbf{W} = \begin{cases} \mathbf{W}_{cr} := \text{Current reputation} \\ \mathbf{W}_{hmm} := \text{HMM-W2V reputation} \\ \mathbf{W}_p := \text{Predictive features in} \\ \quad \text{kokkodis et al. (2015)} \\ \mathbf{W}_{hmm \wedge p} := \mathbf{W}_p \wedge \text{HMM-W2V reputation} \end{cases} \quad (11)$$

Online Appendix G presents the list of predictive features \mathbf{W}_p .

I evaluate the performance of each recommender by estimating their area under the curve (AUC). To highlight the benefits of using HMM-W2V reputation, I estimate the following improvements:

$$\begin{aligned} &\text{Improvement over current reputation} \\ &= \frac{AUC(\mathbf{W}_{hmm}) - AUC(\mathbf{W}_{cr})}{AUC(\mathbf{W}_{cr})} * 100, \\ &\text{Improvement over predictive features} \\ &= \frac{AUC(\mathbf{W}_{hmm \wedge p}) - AUC(\mathbf{W}_p)}{AUC(\mathbf{W}_p)} * 100. \end{aligned}$$

Figure 9 shows the 10-fold cross-validated AUC improvements. First, compared with the current reputation, HMM-W2V reputation provides significantly ($p < 0.001$) better recommendations across all classifiers that yield an AUC improvement between 2.4% and 10%. Second, once I include the HMM-W2V reputation in the set of predictive features presented in Online Appendix G, the performance of the recommenders increases between 1.3% ($p < 0.001$) for neural

networks and 3.5% ($p < 0.001$) for logistic regression. Such AUC improvements could generate better matches, happier employers, better collaborations, and a subsequent increase in market revenue (Kokkodis et al. 2015).⁴

To conclude, the proposed reputation framework outperforms adaptations of recommender systems, further highlighting the benefits of using the three-component architecture that addresses reputation attribution, reputation staticity, and reputation inflation. Furthermore, Figure 9 illustrates that reputation and recommender systems are not at odds; instead, they provide better user experience when they work together.

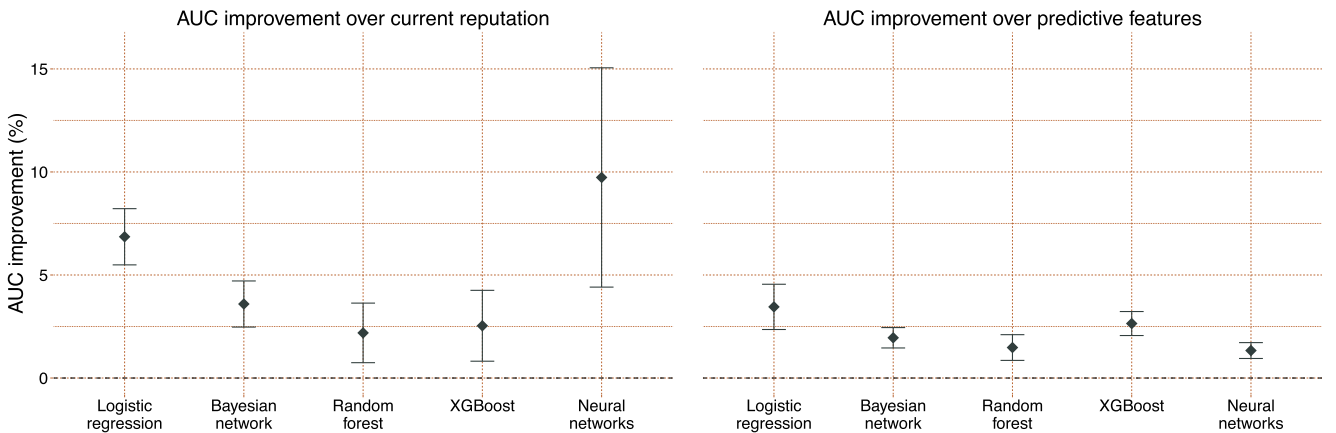
5.5. Additional Evaluations and Generalizability

Besides this analysis, I further evaluate the predictive and explanatory performance of the HMM-W2V framework and its generalizability.

Predictive and Explanatory Performance. Online Appendix D provides performance evaluations and comparisons across different metrics. Specifically, Figure 13 shows that the HMM-W2V framework significantly ($p < 0.001$) outperforms all alternative reputation systems in terms of predictive mean absolute error (MAE) and root mean squared error (RMSE). Even further, Table 5 shows that, compared with alternative reputation systems, HMM-W2V reputation better explains the variance of the observed performance in terms of R^2 in a linear regression specification.

Generalizability. One advantage of the focal approach is that it can generalize to other contexts that experience reputation attribution, reputation staticity,

Figure 9. (Color online) HMM-W2V Reputation Enhances the Performance of Job-Applicant Recommender Systems



Notes. Using HMM-W2V reputation in recommender systems that rank job applicants according to their likelihood of getting hired (Kokkodis et al. 2015, Abhinav et al. 2017) results in better ($p < 0.001$) recommendations. The y-axis shows the AUC improvement when each recommender system on the x-axis includes HMM-W2V reputation as an attribute. This example illustrates how reputation and recommender systems can work together to improve transaction efficiency in a marketplace. The figure shows 10-fold cross-validation average improvements. Error bars represent 95% confidence intervals.

and reputation inflation. One such example is online reputation platforms (e.g., TripAdvisor, Yelp). These platforms deliver positively skewed reputation scores (Luca and Zervas 2016). At the same time, because venues evolve (e.g., change their menus, go through renovations, hire different personnel), their reputation is also dynamic. Finally, these restaurants receive evaluations for multiple dimensions, and, as a result, their reputation systems likely experience reputation attribution as well. Online Appendix E.3 implements the HMM-W2V framework and the alternative reputation and recommender systems on 77,044 restaurant reviews from a major restaurant review platform: The HMM-W2V framework provides significantly ($p < 0.001$) better restaurant reputation scores compared with alternative reputation systems (Figure 14) and adaptations of popular recommender systems (Figure 15). These results empirically show the generalizability of the HMM-W2V framework.

One concern of the proposed approach is that platforms are already developing multidimensional reputation systems, and, as a result, the HMM-W2V framework might be only recovering noisy information from such human-rated dimensions. The focal restaurant review platform creates a perfect environment to test this, as it allows customers to rate venues across four dimensions: “Food,” “Atmosphere,” “Value,” and “Service.” Online Appendix E.3 and Figure 16 compare the HMM-W2V framework with alternative approaches that use as input these human-rated reputation dimensions. The HMM-W2V framework significantly ($p < 0.001$) outperforms these alternatives, suggesting that the latent dimensions captured by the HMM-W2V framework contain different information than the observed multidimensional reputation scores.

Overall, the empirical analysis in this section shows the advantages of the HMM-W2V reputation over a set of advanced alternative reputation and recommender systems in terms of ranking workers, identifying “nonperfect” workers, generating better within-opening rankings, and generating a more representative, normal-like reputation distribution.

6. Discussion

Current reputation systems in online labor markets experience three shortcomings: (1) reputation attribution (attribution of reputation to specific skills is infeasible), (2) reputation staticity (the assumption that worker quality is static and does not evolve), and (3) reputation inflation (feedback scores are positively skewed). To address these shortcomings, this work presents the HMM-W2V framework, which combines human input (assigned feedback scores) with machine learning (word embedding, hidden Markov models) to provide accurate quality estimates for any online worker on any given set of skills. The proposed

framework includes three components. The first component maps skills into a latent space of finite competency dimensions and addresses reputation attribution. The second builds dynamic competency-specific quality assessment models and addresses reputation staticity. The final component aggregates these competency-specific assessments to generate a reputation score for any given set of skills. Because these reputation scores are skillset-specific and evolve over time, they generate a representative, closer-to-normal reputation distribution, and hence address reputation inflation. Application of the proposed framework to two different data sets illustrates that, compared with alternative reputation systems, HMM-W2V reputation performs better in terms of (1) ranking workers according to their likelihood of performing well, (2) identifying “nonperfect” workers who are more likely to underperform and are harder to predict, (3) improving the ranking of within-opening choices, and (4) creating closer-to-normal reputation distributions that facilitate worker differentiation.

6.1. Research Contributions

Given the projected growth of the number of online workers in the coming years (Agile-1 2016, Sundararajan 2016) and their dynamic nature (Oliver 2015, Kokkodis and Ipeirotis 2016), accurate quality assessment could be a defining factor of the ultimate reach of online work. This paper is the first to outline shortcomings of current reputation systems (i.e., reputation staticity and reputation attribution) and to explain why such systems underperform in this context. By identifying these shortcomings, this work provides a solution that generalizes to arbitrary sets of skills. Because it makes skillset-specific estimates that dynamically evolve, this work provides accurate skillset-specific reputation.

From a design perspective, this work extends the rich literature (Table 1) of reputation systems by combining human input with machine intelligence to enhance employer judgment in choosing appropriate online workers. Compared with previous human-based, machine-based, and hybrid reputation systems, the proposed approach has unique dynamic attributes that fit the particular context of online work. Given that algorithmic collaboration between humans and machines is expected to grow (Jain et al. 2018), the proposed hybrid approach could be a baseline for future intelligence augmentation systems that could potentially use human input beyond the current uniform worker assessment (e.g., by having employers rating workers only in dimensions that employers have appropriate expertise in, or by using expert, third-party human raters).

6.2. Methodological Contributions and Generalizability

Methodologically, this paper provides detailed guidelines for markets that are interested in developing dynamic reputation systems by addressing methodological challenges that include the conceptualization, modeling, and estimation of a worker's reputation. Specifically, these guidelines are the following:

- ◊ *Skillset decomposition*: This paper is the first to conceptualize skillsets as documents and propose the application of a text-analysis algorithm in a completely different context. As discussed in Sections 2 and 3 and illustrated in Online Appendix C and Figure 11, this decomposition is necessary, as it allows the proposed reputation framework to generalize on any number of available skills.

- ◊ *HMM architecture*: Section 3.2 describes how practitioners can conceptualize and formulate a suitable structure for an HMM that allows a series of observed signals to shape the transition and emission probabilities of workers of various qualities.

- ◊ *Parameter estimation*: Online Appendix B guides practitioners through the derivation of the global likelihood of the model and the estimation process of all the parameters.

- ◊ *Design choices and evaluation*: Appendix C presents the process of evaluating different components of the HMM-W2V framework and selecting an appropriate configuration of the HMM. Section 5 and Online Appendices C and D guide practitioners on how to evaluate and compare alternative design choices and reputation systems.

These methodological contributions generalize beyond the focal context of online work. The HMM-W2V framework can be adjusted and implemented in any online platform that experiences reputation attribution, reputation staticity, and reputation inflation. For instance, Online Appendix E.3 shows that, under the assumption that hotels and restaurants dynamically change, reputation platforms such as Yelp and TripAdvisor could use a similar framework to develop a more dynamic, up-to-date reputation system. Sharing economy platforms such as Uber and Lyft can borrow ideas from the proposed approach and develop similar reputation systems internally. Such internal systems could estimate the dynamic service quality of each driver and even identify heterogeneity in reputation across various types of trips (e.g., airport and train station trips, long trips out of town, trips in rush hour, Saturday night trips). Finally, online platforms that track workers' career paths, such as LinkedIn, can adapt the proposed approach and estimate the evolving expertise of their users across multiple skillsets. Recommender systems could then use such estimates to suggest new skills or promote workers to potentially relevant jobs.

6.3. Implications for Platforms, Workers, Employers, and the Future of Work

Online labor platforms stand to benefit through implementing the proposed approach, as accurate reputation scores (1) help workers to differentiate, (2) guide employers to make informed and fast (reduced search cost; see Bakos 1997) decisions, and (3) enable the market to improve its recommendation algorithms. When workers can be accurately differentiated, the quality of the supply side of the market naturally increases. High-quality relevant workers are more likely to keep participating, as they are in high demand. Lower-quality workers could be motivated to invest in different skills that are uncorrelated with their current skillset and respective reputation. At the same time, and as I showed in Section 2.3 and Figure 9, markets can improve the performance of their recommendation algorithms by using HMM-W2V reputation as an additional attribute. Better recommendations imply higher income and higher transaction efficacy (Kokkodis et al. 2015).

The performance of the proposed approach in terms of identifying "nonperfect" workers is of particular importance to market managers. Accurately predicting underperformance allows informing employers preemptively. Such interventions could potentially reduce the number of adverse outcomes. Employers who make better-informed and faster decisions that lead to better outcomes are more likely to be happy and keep participating in the marketplace, thereby generating a continuous stream of revenue for the platform (Tripp and Grégoire 2011).

Through the proposed reputation framework, platforms can better understand the supply distributions across latent competencies and any arbitrary combination of skills. Based on such information, market managers can intervene where they deem appropriate (e.g., through targeted advertising on competencies that workers tend to underperform). Online Appendix H empirically analyzes the focal dimensions and explains how market managers can track worker performance and employer demand across competencies and devise appropriate interventions.

The proposed reputation framework combines human input (feedback scores) with machine learning (AI) to augment intelligence in decision making. Over time, through continuous training, the AI component of the framework will improve its performance. As the framework becomes better and more accurate, its effect on humans (both workers and employers) who interact with it will also intensify. Specifically, the AI framework will be facilitating increased differentiation of intelligent crowd and crowd abilities. Due to better differentiation, employers who interact with the system will be able to make intelligent decisions that likely lead to successful outcomes. Such outcomes could encourage

participation in these markets and attract new employers. Similarly, workers that the system evaluates to have high relevant expertise will be able to find tasks to complete online seamlessly. On the other hand, the framework might marginalize workers that it does not evaluate as experts. For such workers, career development systems can provide recommendations for new skills to learn (Kokkodis and Ipeirotis 2020). At the same time, as the IA framework's performance improves, it will potentially allow new workers (without prior history on the platform) to complete tasks for which the system deems workers with history on the platform as inadequate. As a result, a better IA system will likely help workers with high-in-demand skills and abilities to succeed, and it will potentially drive workers who exercise low-demand skills to consider reskilling.

These effects of augmented intelligence could extend to offline work, as the proposed framework could generalize (for instance, through LinkedIn) to offline worker reputation. As automation keeps evolving, the nature of many jobs will change, and other jobs will become obsolete (Brynjolfsson et al. 2018). This transition will require many workers to learn new skills: By some estimates, 120 million workers worldwide will need to be retrained as a result of automation in the next three years (IBM Institute of Business Value 2019). Given this dynamic evolution of skills and workers, a reputation framework such as the one proposed in this work could successfully facilitate supply redistribution while intelligently differentiating relevant hireable workers.

6.4. Additional Discussion of the Proposed Framework

To provide dynamic recommendations, the proposed approach assumes that the competency-specific hidden quality states are discrete. In particular, the HMMs allow workers to transition across these discrete states as they complete new tasks. Once an HMM estimates the state of a worker, it predicts quality estimates that are continuous through Equation (6) and subsequently Equation (8). Hence, although the HMM-W2V framework assumes discrete hidden states, it provides scores that are continuous and capture the complete spectrum of worker competency-specific qualities. Future work could also explore state-space models (e.g., linear dynamical systems) that allow hidden states to be continuous (Murphy 2012). Given that the HMM-W2V framework already provides continuous worker scores, the additional complexity of these models does not guarantee better performance. Even further, latent-space models might not provide clear market insights to managers: Because workers will not reside in discrete latent states, competency-specific market analysis (such as the one in Online Appendix H) will require threshold tuning.

Through the examination of different transition functions (Online Appendix A), I concluded that, for the focal data set, multinomial transitions yield better results. By definition, multinomial transitions allow workers to downgrade to lower-quality states. This might appear irrational at first: Why would a worker's quality decrease over time? The way that an HMM works might provide some rationale. Because each HMM focuses on accurately capturing a worker's latent quality, more worker observations will generate less-noisy estimates of the worker's quality. For instance, consider a worker who receives high feedback scores in early tasks and low feedback scores in later tasks. The HMM will (likely) initially assign this worker to a high-quality state. As the HMM evaluates new evidence, it will (likely) adjust the worker to a lower-quality state. As a result, allowing transitions to lower-quality states is vital for the HMM-W2V framework in order to correct and adjust worker-quality estimates as new information arrives. Online Appendix I discusses this process in greater detail, and Figure 18 empirically shows the benefits of facilitating unconstrained HMM transitions.

An additional concern for the proposed approach is that platforms are already developing multidimensional reputation systems, and, as a result, the HMM-W2V framework might be only recovering noisy information from such human-rated dimensions. To investigate, Online Appendix H presents the skills that define each competency in the focal data set. The complexity of these competency-specific skillsets demonstrates the mental load that humans would need to follow to decompose evaluations across different dimensions. Furthermore, in a human-rating scenario, the mapping of each skill to a nonprimary dimension (that the proposed framework does automatically) would be practically infeasible. For instance, if a worker completes a data-mining job, a human rater would evaluate the worker's performance in that skill, but not in every other skill available on the market (e.g., video production). The HMM-W2V framework does this mapping automatically, as it implicitly estimates correlations between skills through the W2V decomposition (Section 3.1). Online Appendix H discusses this in detail, and Online Appendix E.3 empirically shows that the proposed mapping captures different information than the already implemented human-rated dimensions.

Finally, the fact that, compared to D2V, W2V was a more appropriate decomposition approach in the worker-reputation context (Online Appendix C.1) but less appropriate in the restaurant-reputation context (Online Appendix E.3) raises a question of when markets should prefer W2V over D2V. Conceptually, W2V is likely to perform better when the pool of terms is predefined, well-structured, and where every term

includes crucial information. In these scenarios, where each individual term (e.g., a skill) contains crucial information, summing up mappings through Equation (1) should generate more informative representations than D2V mappings that will unavoidably include noise from rare combinations of terms (e.g., rare combinations of skills). On the other hand, D2V should perform better when the analyzed text is unstructured (e.g., review text). In those cases, summing up weights of seemingly random terms (e.g., found in reviews) through Equation (1) will likely generate noisy representations.

6.5. Conclusion

Conclusively, this work presents an intelligence augmentation framework that addresses reputation attribution, reputation staticity, and reputation inflation. Application of this framework in two different contexts (online labor markets and online reputation platforms) shows that it can track evolving entities across multiple dimensions and provide accurate service-quality estimates. As a result, its deployment on different types of online platforms could have significant implications for workers, employers, businesses, and the future of work.

Acknowledgments

The author thanks Sam Ransbotham, Panagiotis Adamopoulos, and Vilma Todri for their guidance and feedback. The author also thanks Robert Fichman, Gerald Kane, Zhuoxin Li, Xuan Ye, and Mike Teodorescu for their comments and suggestions on improving the paper.

Endnotes

¹ Allahbakhsh et al. (2012) weight the timing of each rating, so, in theory, it has the potential to address reputation staticity. Table 1 clarifies this point.

² Alternatively, workers could land stochastically to any of the available states. This would add noise to the estimates, and as a result it could potentially hurt the performance of the framework.

³ Vectors \mathbf{Z}_{t-1} and \mathbf{X}_t (which I use later) are worker-specific. As a result, the transition and emission probabilities are also worker-specific. For simplicity, I drop worker subscript i from the notation of this subsection. I use the subscript i in Section 3.3 and Online Appendix B to highlight that the skillset-specific reputation scores are also worker-specific.

⁴ The high variance of the neural networks performance when modeling only reputation scores (i.e., only one feature) suggests that the networks possibly overfit the training sets and often fail in the test sets. Once additional features enter the model, trained neural networks make robust predictions that perform better than alternative approaches.

References

Abhinav K, Dubey A, Jain S, Virdi G, Kass A, Mehta M (2017) Crowdadvisor: A framework for freelancer assessment in online marketplace. *IEEE/ACM 39th Internat. Conf. Software Engrg.: Software Engrg. Practice Track* (IEEE, New York), 93–102.

Adamopoulos P (2013) Beyond rating prediction accuracy: on new perspectives in recommender systems. *RecSys '13: Proc. 7th ACM Conf. Recommender Systems* (ACM, New York), 459–462.

Adamopoulos P, Tuzhilin A (2013) Recommendation opportunities: improving item prediction using weighted percentile methods in collaborative filtering systems. *RecSys '13: Proc. 7th ACM Conf. Recommender Systems* (ACM, New York), 351–354.

Adamopoulos P, Tuzhilin A (2014) On over-specialization and concentration bias of recommendations: Probabilistic neighborhood selection in collaborative filtering systems. *RecSys '14: Proc. 8th ACM Conf. Recommender Systems* (ACM, New York), 153–160.

Adamopoulos P, Tuzhilin A (2015) The business value of recommendations: A privacy-preserving econometric analysis. *Proc. 36th Internat. Conf. Inform. Systems* (AIS, Atlanta).

Adomavicius G, Tuzhilin A (2005) Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions. *Trans. Knowledge Data Engrg.* 17(6):734–749.

Agile-1 (2016) The new gig economy: What's at stake? Retrieved December 2, 2019, http://www.hrotoday.com/wp-content/uploads/2016/07/Whitepaper_Agile2016-single.pdf.

Agrawal A, Lacetera N, Lyons E (2013) Does information help or hinder job applicants from less developed countries in online markets? NBER Working Paper 18720, National Bureau of Economic Research, Cambridge, MA.

Agrawal R, Srikant R (1994) Fast algorithms for mining association rules. Bocca JB, Jarke M, Zaniolo CA, eds. *VLDB '94: Proc. 20th Internat. Conf. Very Large Data Bases* (Morgan Kaufmann Publishers, San Francisco), 487–499.

Akerlof GA (1970) The market for “lemons”: Quality uncertainty and the market mechanism. *Quart. J. Econom.* 84(3):488–500.

Allahbakhsh M, Ignjatovic A, Benatallah B, Bertino E, Foo N (2012) Reputation management in crowdsourcing systems. *Internat. Conf. Collaborative Comput.: Networking Appl. Worksharing (CollaborateCom)* (IEEE, New York), 664–671.

Amazon (2018) Customer reviews. Retrieved December 2, 2019, <https://www.amazon.com/gp/help/customer/display.html?nodeId=202094910>.

Amazon (2020) Search by “avg. customer review.” Retrieved May 2, <https://www.amazon.com/gp/help/customer/display.html?nodeId=201889520>.

Autor DH (2001) Wiring the labor market. *J. Econom. Perspect.* 15(1):25–40.

Autor DH, Katz LF, Krueger AB (1998) Computing inequality: Have computers changed the labor market? *Quart. J. Econom.* 113(4): 1169–1213.

Ba S, Pavlou PA (2002) Evidence of the effect of trust building technology in electronic markets: Price premiums and buyer behavior. *MIS Quart.* 26(3):243–268.

Baba Y, Kinoshita K, Kashima H (2016) Participation recommendation system for crowdsourcing contests. *Expert Systems Appl.* 58:174–183.

Bakos Y (1997) Reducing buyer search costs: Implications for electronic marketplaces. *Management Sci.* 43(12):1676–1692.

Balog K, De Rijke M (2007) Determining expert profiles (with an application to expert finding). Sangal R, Mehta H, Bagga RK, eds. *IJCAI '07: Proc. 20th Internat. Joint Conference Artificial Intelligence* (Morgan Kaufmann Publishers, San Francisco), 2657–2662.

Banker RD, Hwang I (2008) Importance of measures of past performance: Empirical evidence on quality of e-service providers. *Contemporary Accounting Res.* 25(2):307–337.

Billsus D, Pazzani MJ (1998) Learning collaborative information filters. Shavlik JW, ed. *ICML '98: Proc. 15th Internat. Conf. Machine Learn.* (Morgan Kaufmann Publishers, San Francisco), 46–54.

Bouguesma M, Dumoulin B, Wang S (2008) Identifying authoritative actors in question-answering forums: The case of Yahoo! answers.

- KDD '08: *Proc. 14th ACM SIGKDD Internat. Conf. Knowledge Discovery Data Mining* (ACM, New York), 866–874.
- Breese SJ, Heckerman D, Kadie C (1998) Empirical analysis of predictive algorithms for collaborative filtering. Cooper GF, Serafin M, eds. *UAI '98: Proc. 14th Conf. Uncertainty Artificial Intelligence* (Morgan Kaufmann Publishers, San Francisco), 43–52.
- Brynjolfsson E, Mitchell T, Rock D (2018) What can machines learn, and what does it mean for occupations and the economy? *AEA Papers Proc.* 108:43–47.
- Chen T, Guestrin C (2016) XGBoost: A scalable tree boosting system. *KDD '16: Proc. 22nd ACM SIGKDD Internat. Conf. Knowledge Discovery Data Mining* (ACM, New York), 785–794.
- Chevalier JA, Mayzlin D (2006) The effect of word of mouth on sales: Online book reviews. *J. Marketing Res.* 43(3):345–354.
- Christoforaki M, Ipeirotis PG (2015) A system for scalable and reliable technical-skill testing in online labor markets. *Comput. Networks* 90:110–120.
- Cui G, Lui H-K, Guo X (2012) The effect of online consumer reviews on new product sales. *Internat. J. Electronic Commerce* 17(1):39–58.
- Curtis N, Safavi-Naini R, Susilo W (2004) X²Rep: Enhanced trust semantics for the XRep protocol. Jakobsson M, Yung M, Zhou J, eds. *Appl. Cryptography Network Security: ACNS 2004* (Springer, Berlin), 205–219.
- Daltayanni M, de Alfaro L, Papadimitriou P (2015) WorkerRank: Using employer implicit judgements to infer worker reputation. *WSDM '15: Proc. 8th ACM Internat. Conf. Web Search Data Mining* (ACM, New York), 263–272.
- Damiani E, di Vimercati DC, Paraboschi S, Samarati P, Violante F (2002) A reputation-based approach for choosing reliable resources in peer-to-peer networks. *CCS '02: Proc. 9th ACM Conf. Comput. Comm. Security* (ACM, New York), 207–216.
- Delgado J, Ishii N (1999) Memory-based weighted majority prediction. *Proc. Special Interest Group Inform. Retrieval Workshop Recommender Systems*.
- Dellarocas C (2003) The digitization of word of mouth: Promise and challenges of online feedback mechanisms. *Management Sci.* 49(10):1407–1424.
- Dellarocas C (2006) Reputation mechanisms. *Handbook on Economics and Information Systems* 629–660.
- Devooght R, Bersini H (2017) Collaborative filtering based on sequences. Retrieved December 2, 2019, <https://github.com/rdevooght/sequence-based-recommendations>.
- Duan W, Gu B, Whinston AB (2008) Do online reviews matter? An empirical investigation of panel data. *Decision Support Systems* 45(4):1007–1016.
- eBay (2018) Feedback forum. Retrieved December 2, 2019, <https://pages.ebay.com/services/forum/feedback.html>.
- Einav L, Farronato C, Levin J (2016) Peer-to-peer markets. *Annual Rev. Econom.* 8:615–635.
- Filippas A, Horton JJ, Golden J (2018) Reputation inflation. *EC '18: Proc. 2018 ACM Conf. Econom. Comput.* (ACM, New York), 483–484.
- Freelancers Union (2017) Freelancing in America. Retrieved December 2, 2019, <https://www.upwork.com/i/freelancing-in-america/2017/>.
- Gandini A, Pais I, Beraldo D (2016) Reputation and trust on online labour markets: The reputation economy of elance. *Work Organ. Labour Globalisation* 10(1):27–43.
- Ghose A, Ipeirotis PG (2011) Estimating the helpfulness and economic impact of product reviews: Mining text and reviewer characteristics. *Trans. Knowledge Data Engrg.* 23(10):1498–1512.
- Ghose A, Ipeirotis PG, Li B (2012) Designing ranking systems for hotels on travel search engines by mining user-generated and crowdsourced content. *Marketing Sci.* 31(3):493–520.
- Ghose A, Ipeirotis PG, Li B (2014) Examining the impact of ranking on consumer behavior and search engine revenue. *Management Sci.* 60(7):1632–1654.
- Goswami A, Hedayati F, Mohapatra P (2014) Recommendation systems for markets with two sided preferences. *Proc. 13th Internat. Conf. Machine Learn. Appl.* (IEEE, New York), 282–287.
- Graham M, Hjorth I, Lehdonvirta V (2017) Digital labour and development: Impacts of global digital labour platforms and the gig economy on worker livelihoods. *Transfer* 23(2):135–162.
- Hambleton RK, Swaminathan H, Rogers HJ (1991) *Fundamentals of Item Response Theory* (Sage Publications, Newbury Park, CA).
- Hendrikx F, Bubendorfer K (2013) Malleable access rights to establish and enable scientific collaboration. *Proc. 9th Internat. Conf. e-Science* (IEEE, New York), 334–341.
- Hendrikx F, Bubendorfer K, Chard R (2015) Reputation systems: A survey and taxonomy. *J. Parallel Distributed Comput.* 75:184–197.
- Ho Y-C, Wu J, Tan Y (2017) Disconfirmation effect on online rating behavior: A structural model. *Inform. Systems Res.* 28(3):626–642.
- Hochreiter S, Schmidhuber J (1997) Long short-term memory. *Neural Comput.* 9(8):1735–1780.
- Horton JJ (2017) The effects of algorithmic labor market recommendations: Evidence from a field experiment. *J. Labor Econom.* 35(2):345–385.
- Hossain MS, Arefin MS (2019) Development of an intelligent job recommender system for freelancers using client's feedback classification and association rule mining techniques. *J. Software* 14(7):312–339.
- Hsueh S-C, Lin M-Y, Chen C-L (2008) Mining negative sequential patterns for e-commerce recommendations. *Proc. Asia-Pacific Services Comput. Conf.* (IEEE, New York), 1213–1218.
- Hu N, Pavlou PA, Zhang J (2017) On self-selection biases in online product reviews. *MIS Quart.* 41(2):449–471.
- Hu N, Zhang J, Pavlou PA (2009) Overcoming the j-shaped distribution of product reviews. *Comm. ACM* 52(10):144–147.
- Huber PJ (2011) *Robust Statistics* (Springer, Berlin).
- IBM Institute for Business Value (2019) The enterprise guide to closing the skills gap. Retrieved December 2, <https://www.ibm.com/downloads/cas/EPYMNBJA>.
- Ipeirotis P (2013) Badges and the Lake Wobegon effect. Retrieved December 30, 2018, <https://www.behind-the-enemy-lines.com/2013/10/badges-and-lake-wobegon-effect.html>.
- Jagabathula S, Subramanian L, Venkataraman A (2014) Reputation-based worker filtering in crowdsourcing. Ghahramani Z, Welling M, Cortes C, Lawrence N, Weinberger KQ, eds. *Advances in Neural Information Processing Systems* (Curran Associates, Red Hook, NY), 2492–2500.
- Jain H, Padmanabhan B, Pavlou PA, Santanam RT (2018) Call for papers—special issue of information systems research—humans, algorithms, and augmented intelligence: The future of work, organizations, and society. *Inform. Systems Res.* 29(1):250–251.
- Jannach D, Lerche L, Jugovac M (2015) Adaptation and evaluation of recommendations for short-term shopping goals. *RecSys '15: Proc. 9th ACM Conf. Recommender Systems* (ACM, New York), 211–218.
- Jerath K, Fader PS, Hardie GSB (2011) New perspectives on customer “death” using a generalization of the pareto/nbd model. *Marketing Sci.* 30(5):866–880.
- Jøsang A, Golbeck J (2009) Challenges for robust trust and reputation systems. *Proc. 5th Internat. Workshop Security Trust Management*.
- Jøsang A, Guo G, Pini MS, Santini F, Xu Y (2013) Combining recommender and reputation systems to produce better online advice. Torra V, Narukawa Y, Navarro-Arribas G, Megías D, eds. *Proc. Internat. Conf. Modeling Decisions Artificial Intelligence* (Springer, Berlin), 126–138.
- Jøsang A, Ismail R, Boyd C (2007) A survey of trust and reputation systems for online service provision. *Decision Support Systems* 43(2):618–644.
- Jurczyk P, Agichtein E (2007) Discovering authorities in question answer communities by using link analysis. *CIKM '07: Proc. 16th*

- ACM Conf. Inform. Knowledge Management (ACM, New York), 919–922.
- Kamvar SD, Schlosser MT, Garcia-Molina H (2003) The EigenTrust algorithm for reputation management in P2P networks. WWW '03: Proc. 12th Internat. Conf. World Wide Web (ACM, New York), 640–651.
- Kanat I, Hong Y, Santanam RT (2018) Surviving in global online labor markets for IT services: A geo-economic analysis. Inform. Systems Res. 29(4):893–909.
- Kantor BP, Rokach L, Ricci F, Shapira B, eds. (2011) *Recommender Systems Handbook* (Springer, Boston).
- Kendall MG (1938) A new measure of rank correlation. Biometrika 30(1–2):81–93.
- Kohl M (2019) Totalvardist: Generic function for the computation of the total variation distance of two distributions. Retrieved December 2, <https://www.rdocumentation.org/packages/distrEx/versions/2.5/topics/TotalVarDist>.
- Kokkodis M (2018) Dynamic recommendations for sequential hiring decisions in online labor markets. KDD '18: Proc. 24th ACM SIGKDD Internat. Conf. Knowledge Discovery Data Mining (ACM, New York), 453–461.
- Kokkodis M (2020) Diversify or specialize? Skillset diversification in online labor markets: Reputation losses and opportunity gains. Working paper.
- Kokkodis M, Ipeirotis PG (2014) The utility of skills in online labor markets. Proc. 35th Internat. Conf. Inform. Systems.
- Kokkodis M, Ipeirotis PG (2016) Reputation transferability in online labor markets. Management Sci. 62(6):1687–1706.
- Kokkodis M, Ipeirotis PG (2020) Demand-aware career path recommendations: A reinforcement learning approach. Manage. Sci. Forthcoming.
- Kokkodis M, Lappas T (2020) Your hometown matters: Popularity-difference bias in online reputation platforms. Inform. Systems Res. 31(2):412–430.
- Kokkodis M, Lappas T, Kane G (2020a) Direct and indirect effects of introducing purchase verification in e-commerce platforms. Working paper.
- Kokkodis M, Lappas T, Ransbotham S (2020b) From lurkers to workers: Predicting voluntary contribution and community welfare. Inform. Systems Res. 31(2):607–626.
- Kokkodis M, Papadimitriou P, Ipeirotis PG (2015) Hiring behavior models for online labor markets. WSDM '15: Proc. 8th ACM Internat. Conf. Web Search Data Mining (ACM, New York), 223–232.
- Kokkodis M, Ransbotham S (2020) Asymmetric reputation spillover from agencies on digital platforms. Working paper.
- Koren Y (2010) Collaborative filtering with temporal dynamics. Commun. ACM 53(4):89–97.
- Koren Y, Bell R, Volinsky C (2009) Matrix factorization techniques for recommender systems. Computer 42(8):30–37.
- Kuhn P, Skuterud M (2004) Internet job search and unemployment durations. Amer. Econom. Rev. 94(1):218–232.
- Kula M (2018) Deep recommender models using PyTorch. Retrieved December 2, 2019, <https://github.com/maciejkula/spotlight>.
- Kung-Hsiang H (2018) Introduction to recommender systems: Part 2 (neural network approach). Towards Data Sci., February 16, <https://towardsdatascience.com/introduction-to-recommender-system-part-2-adoption-of-neural-network-831972c4cbf7>.
- Le Q, Mikolov T (2014) Distributed representations of sentences and documents. Proc. 31st Internat. Conf. Machine Learning (JMLR, Cambridge, MA), 1188–1196.
- Lee G, Raghu TS (2014) Determinants of mobile apps' success: Evidence from the appstore market. J. Management Inform. Systems 31(2):133–170.
- Lemire D, Maclachlan A (2005) Slope one predictors for online rating-based collaborative filtering. Proc. Internat. Conf. Data Mining (SIAM, Philadelphia), 471–475.
- Lin M, Liu Y, Viswanathan S (2016) Effectiveness of reputation in contracting for customized production: Evidence from online labor markets. Management Sci. 64(1):345–359.
- Longadge R, Dongre S (2013) Class imbalance problem in data mining review. Preprint, submitted May 8, <https://arxiv.org/abs/1305.1707>.
- Lu X, Ba S, Huang L, Feng Y (2013) Promotional marketing or word-of-mouth? Evidence from online restaurant reviews. Inform. Systems Res. 24(3):596–612.
- Luca M (2016) Reviews, reputation, and revenue: The case of yelp.com. Working paper, Harvard University, Boston.
- Luca M (2017) Designing online marketplaces: Trust and reputation mechanisms. Innovation Policy Econom. 17:77–93.
- Luca M, Zervas G (2016) Fake it till you make it: Reputation, competition, and Yelp review fraud. Management Sci. 62(12):3412–3427.
- Meyer F (2012) Recommender systems in industrial contexts. Preprint, submitted March 20, <https://arxiv.org/abs/1203.4487>.
- Mikolov T, Chen K, Corrado G, Dean J (2013) Efficient estimation of word representations in vector space. Preprint, submitted January 16, <https://arxiv.org/abs/1301.3781>.
- Moe WW, Schweidel DA (2012) Online product opinions: Incidence, evaluation, and evolution. Marketing Sci. 31(3):372–386.
- Moreno A, Terwiesch C (2014) Doing business with strangers: Reputation in online service marketplaces. Inform. Systems Res. 25(4):865–886.
- Murphy KP (2012) *Machine learning: A Probabilistic Perspective* (MIT Press, Cambridge, MA).
- Nica E, Potcovaru A-M, Mirică C-O (2017) A question of trust: Cognitive capitalism, digital reputation economy, and online labor markets. Econom. Management Financial Markets 12(3):64–69.
- Oliver B (2015) Redefining graduate employability and work-integrated learning: Proposals for effective higher education in disrupted economies. J. Teaching Learning Graduate Employability 6(1):56–65.
- Pallais A (2014) Inefficient hiring in entry-level labor markets. Amer. Econom. Rev. 104(11):3565–3599.
- Paolacci G, Chandler J, Ipeirotis PG (2010) Running experiments on amazon mechanical turk. Judgment Decision Making 5(5):411–419.
- Patel B, Kakuste V, Eirinaki M (2017) CaPaR: A career path recommendation framework. Proc. Internat. Conf. Big Data Comput. Service Appl. (BigDataService) (IEEE, New York), 23–30.
- Pavlou PA, Gefen D (2004) Building effective online marketplaces with institution-based trust. Inform. Systems Res. 15(1):37–59.
- Pelechrinis K, Zadorozhny V, Kounev V, Oleshchuk V, Anwar M, Lin Y (2015) Automatic evaluation of information provider reliability and expertise. World Wide Web 18:33–72.
- Proserpio D, Zervas G (2017) Online reputation management: Estimating the impact of management responses on consumer reviews. Marketing Sci. 36(5):645–665.
- Quadrana M, Cremonesi P, Jannach D (2018) Sequence-aware recommender systems. ACM Comput. Surveys 51(4):66.
- Rahman H (2018a) Don't worship the stars: Ratings inflation in online labor markets. Proc. Internat. Conf. Inform. Systems (AIS, Atlanta).
- Rahman HA (2018b) Reputational ploys: Reputation and ratings in online markets. Acad. Management Proc. 2018(1).
- Rendle S, Freudenthaler C, Schmidt-Thieme L (2010) Factorizing personalized Markov chains for next-basket recommendation. WWW '10: Proc. 19th Internat. Conf. World Wide Web (ACM, New York), 811–820.
- Ricci F, Rokach L, Shapira B (2011) Introduction to *Recommender Systems Handbook*. Ricci F, Rokach L, Shapira S, Kantor PB, eds. *Recommender Systems Handbook* (Springer, Boston), 1–35.
- Sabater J, Sierra C (2001a) REGRET: Reputation in gregarious societies. AGENTS '01: Proc. 5th Internat. Conf. Autonomous Agents (ACM, New York), 194–195.

- Sabater J, Sierra C (2001b) Social ReGrE, a reputation model based on social relations. *ACM SIGecom Exchanges* 3(1):44–56.
- Schmidt FL, Hunter JE (1983) Individual differences in productivity: An empirical test of estimates derived from studies of selection procedure utility. *J. Appl. Psych.* 68(3):407–414.
- Smola AJ, Schölkopf B (2004) A tutorial on support vector regression. *Statistics Comput.* 14:199–222.
- Spearman C (1904) The proof and measurement of association between two things. *Amer. J. Psych.* 15:72–101.
- Stack Overflow (2018) What is reputation? How do I earn (and lose) it? Retrieved December 2, 2019, <https://stackoverflow.com/help/whats-reputation>.
- Stevenson B (2009) The Internet and job search. Autor DH, ed. *Studies of Labor Market Intermediation* (University of Chicago Press, Chicago), 67–86.
- Su X, Khoshgoftaar TM (2009) A survey of collaborative filtering techniques. *Adv. Artificial Intelligence* 2009:4.
- Sundararajan A (2016) *The Sharing Economy: The End of Employment and the Rise of Crowd-Based Capitalism* (MIT Press, Cambridge, MA).
- Tadelis S (2016) Reputation and feedback systems in online platform markets. *Annu. Rev. Econom.* 8:321–340.
- Tian C, Yang B (2011) R2Trust, a reputation and risk based trust management framework for large-scale, fully decentralized overlay networks. *Future Generation Comput. Systems* 27(8): 1135–1141.
- Tripp MT, Grégoire Y (2011) When unhappy customers strike back on the Internet. *MIT Sloan Management Rev.* 52:37–44.
- Turkopticon (2018) Turkopticon. Retrieved December 2, 2019, <https://turkopticon.ucsd.edu/>.
- Upwork (2014) Online work report. Retrieved April 28, 2019, <https://web.archive.org/web/20180228011632/http://elance-odesk.com:80/online-work-report-global>.
- Upwork (2018) Add certifications. Retrieved December 2, 2019, <https://support.upwork.com/hc/en-us/articles/215650138-Add-Certifications>.
- Wood-Doughty A (2018) The role of reputation systems in an online labor market. Working paper.
- Xie H, Lui JCS, Towsley D (2015) Incentive and reputation mechanisms for online crowdsourcing systems. *Proc. 23rd Internat. Sympos. Quality Service (IWQoS)* (IEEE, New York), 207–212.
- Xiong L, Liu L (2004) PeerTrust: Supporting reputation-based trust for peer-to-peer electronic communities. *Trans. Knowledge Data Engrg.* 16(7):843–857.
- Ye Q, Law R, Gu B (2009) The impact of online user reviews on hotel room sales. *Internat. J. Hospital Management* 28(1):180–182.
- Yoganarasimhan H (2013) The value of reputation in an online freelance marketplace. *Marketing Sci.* 32(6):860–891.
- Zang H, Xu Y, Li Y (2010) Non-redundant sequential association rule mining and application in recommender systems. *Proc. IEEE/WIC/ACM Internat. Conf. Web Intelligence Intelligent Agent Tech.*, vol.3. (IEEE, New York), 292–295.
- Zervas G, Proserpio D, Byers J (2015) A first look at online reputation on Airbnb, where every stay is above average. Working paper.
- Zhang J, Ackerman MS, Adamic L (2007) Expertise networks in online communities: structure and algorithms. *WWW '07: Proc. 16th Internat. Conf. World Wide Web* (ACM, New York), 221–230.